**Advanced Molecular Dynamics Simulation Techniques for Kinetic Analysis of Biological Systems**

**Wouter Vervust**

GHENT
UNIVERSITY

# Members of the Examination Board

# ACKNOWLEDGMENTS

My PhD journey started in August 2020, about four years and one month ago, or to frame it within a more memorable timeline: somewhere in the middle of the 'inter-lockdown' COVID period in Belgium. On the first day, I was appointed a spot in the office of Ghazal, Carlos, and Samaneh, where the main side-quest of a PhD student dawned upon me—as upon many before me; to outlast your colleagues and claim their superior desks. Over the years I became quite proficient at *desk-hopping*, until three years into the PhD my migratory existence could turn sedentary in Carlos' glorious spot. A short-lived victory, as quickly thereafter our BioMMedA group moved to *The Core* and later to *The Coretainers*, where us BioMMediAns currently reside. This is far from a complaint, as I often liked to work from different locations, even outside of office spaces. Going one step further, I even liked these changes in scenery, as they provide a natural moment to reflect on the past, and by doing so, provide a new (and hopefully improved) perspective. But enough office talk—after all, this is my final out-of-office notice, and I'd like to thank the people that have supported me on this journey.

Prof. An Ghyels, half-way through my Master's thesis you had brought the news of switching your research focus to biological systems. Retrospectively, this was not as surprising of news, as it confirms a quality that I had already encountered (and deeply respected) in you, being your intellectual curiosity. I required little contemplation to apply for the open PhD position in your newly formed lab, as prior long ponderings had already revealed that I wanted to learn new things, to be a PhD-*student* in a new field. For this I want to express my deep gratitude to you, An, for providing me with this pursuit of shared curiosity, in which you have supported me through the exploration of new ideas. You were my only supervisor, and this definitely did not equate in a lack of research ideas, nor in a lack of supervision. In contrast, your long list of research ideas has rendered

this my PhD to be somewhat dual in nature, a fact for which I am very appreciative (more to learn!).

Next, I want to express my gratitude to the doctoral exam committee members, prof. Gert De Cooman, prof. Titus van Erp, prof. Jocelyne Vreede, prof. Louis Vanduyfhuys, and prof. Paul Van Liedekerke. Thank you for your time and effort in examining this dissertation, and for the interesting discussions that arose during the internal defense.

Titus, I am especially grateful to have had such close collaboration with you. Your ability to explain complex technical and theoretical concepts in an easily digestible manner has helped enormously, and this work would not have reached this level without you. I continue to look up to you, and being on-par with your jokes is definitely my next milesto-, euh, interface ;).

I want to thank the other wonderful path-samplers I've had the pleasure to meet during my stays at the NTNU: Enrico, Daniel, and Lukas. Enrico, without your aid the implementation of partial paths in PyRETIS would not have been possible. I hope that one day we meet again, on or off a race bike. Daniel, I hope you still remember my name with all those Vim keybindings in your head. Thank you for all the coding help, that monstrosity 'the july-branch', and inviting us to your house. When I revisit Japan, I'll try to pass by! Lukas, I hope you have had some time in-between all the other projects to work on QuanTIS. It was definitely worth to travel out at night to watch that weak Aurora, or was that just one of my dreams?

Gert, it was a pleasure to assist in SySi together with Arne, Keano and Martijn. I had a lot of fun during our weekly chess-puzzle sessions (and the exercise lectures that accompanied it), and the lunches/dinners. I'll make sure to create a fresh batch of Jean-Luc Picard questions for the generations to come.

I want to thank the colleagues of our small molecular modeling group. Samaneh and Bhawna, together we formed the 'first-generation' molecular modelers at BioMMedA, and it has been a pleasure. Samaneh, I now realize that our first and latest conversation was about you being 24, the age that you shall never surpass. I hope you are enjoying your post-doc in Switzerland, and that have befriended the local Water. Bhawna, you have as many ideas going through your head as you have Chrome tabs open, making conversations with you interesting and unpredictable. Thanks for always being brutally honest and spontaneous. Sina and Elias, although we only

met during the final year of my PhD, it felt good to have other people directly working on path sampling. Sina, thanks for all the interesting talks during our on-one-on lunches. I appreciated your technical computing skills, your laid-backness, and your dry-and-direct humor. Elias, I appreciated the many interesting path sampling conversations, and the bugs you caught in the toytis code. You shared my somewhat weird sense of humor, resulting in some hilarious moments. Remember to stay away from the Green Ones, and to invite me to your CRIG presentation about staple-paths.

Next, I want to thank my colleagues of BioMMedA, the research group in which we were embedded, and our neighboring group Medisip.

Yousof and Amith, I quickly realized that we could delve into nearly any topic without it ever becoming dull. Naturally, this led to us spending more time together outside of work, where I came to see you both as friends rather than colleagues. Yousof, I admired in particular your forthrightness, regardless of the circumstances, where your candid nature has had a significant impact on the way I view many things. I'd also like to extend my gratitude to your father, Mohamed Assaad, who is *by far* the most hospitable person I have ever met, which made our stay in Cairo all the more memorable. I'm certain another adventure together is just around the corner. Amith, my lawfully wedded wife, your knack for playing devil's advocate never ceased to impress me. Your satirical, witty, self-deprecating, and sometimes outrageous jokes were always music to my ears. While you often claim to be reserved, I have at many times witnessed your ability 'to strike a conversation with just about anyone', albeit occasionally accompanied by an ulterior motive ;).

I want to thank Jens, Meysam and Stefaan for the regular UZ resto visits, which offered a nice variety of serious and silly conversations. Jens, thanks for maintaining my love for physics. Meysam, thanks for radiating happiness, and congratulations on 'your' epic sense of fashion. Stefaan, thank you for being the cool professor that you are, and for letting me use your race bike for nearly two years. Maya and Rabia, you often joined to the resto, where your presence was always enjoyable. Melissa, your unique sense of humor made the journey all the more enjoyable, and I think the time has come for you to add another washable title on my Wiener Melange cup. Also thanks to the other Medisip members, Jolan, Amir, Boris, Emma, and Florance for our short albeit nice crossings.

Sarah, thank you for having no filter whatsoever, and teaching me how to your 'marginale self' take control. Or was it the other way around? Hooman and Mohammad, I often tried to infiltrate your light-hearted 3pm 'Persian' breaks with Samaneh, Amir, and Sina. I now claim that the writing of the 'Digesting duck' is a direct result of listening to Shahjarian, father or son. Tim, the unofficial manager and representative of BioMMedA, thank you for imposing your good mood on the office and the social gatherings you so often planned. Simeon, the guardian and muscle power of BioMMedA, thank you for all the historical facts, and your numerous contributions to the WoTD. Saar, your silly humor was much appreciated, may you forever flex on us mortals with your hundreds of gold medals. Lise, thank you for being as organized as you are, and thereby inadvertently shifting me slightly away from chaos. Matthias, thanks for always providing the half-full angle to many glasses. Ariana, when later in life you get to host your own *loft* parties, don't forget your roots and invite me. Jessie, consistently arriving at 8 (or was it 8.30? I'm not too familiar with these realms) is a very impressive feat, and thanks for helping out with SySi on such short notice. Jellis, 'the fastest in the West'. Although it's limited to the context of periodically swinging your legs without breaking contact with the ground, it's a nice title to hold.

Successfully finishing my PhD would not have been possible without the foundation of a full life provided by a special number of people. I want to specifically thank my dearest group of friends: David, Maarten, Pieter-Jan, Gilles, Gert, Dario, Wouter and Daan. Without you, I cannot even imagine *who* I would be. It has been beautiful to see how we evolved, both as individuals and as a group, and it will be especially beautiful to see where this evolution brings us. Antonino, I want to thank you specifically for fostering my interest in philosophy. You are among the most creative and spontaneous people I have met, and I'm *sure* that both you and I will enjoy this bond much longer than you think.

Graag wil ik ook mijn familie, de Vosjes, bedanken: Carl voor de serieuze gesprekken, en de, jah, niet zo-serieuze gesprekken ;), Sylvie voor de nuchtere kijk op alles, en Marie, Irma, en Elodie voor het familie-leven uit te breiden tot een vrienden-leven, waar groot Flavietje later ook wel deel van zal uitmaken!

Marraine en Peter, jullie zijn altijd dicht betrokken geweest op mijn leven, beide letterlijk en figuurlijk. Marraine, ik wil je speciaal bedanken voor de waarden die ik dankzij jou hoog in het vaandel draag. Ik ken alvast één iemand die het niet graag zal horen, maar

jouw frietjes en pompoensoep zijn de beste die er zijn ;). Peter, jouw ingenieuze creativiteit en jouw kunde om je mannetje te staan in deze wereld zijn grote inspiraties geweest in mijn leven. Tot op vandaag blijf je een groot voorbeeld waar ik naar op kijk.

Mama en papa, zonder jullie geen 'ik', en die 'ik' die jullie met zoveel zorg gevormd hebben, is jullie dankbaar op een manier die moeilijk te verwoorden is. Doorheen mijn leven hebben jullie mij onvoorwaardelijk gesteund, waar het doctoraatsgedeelte dan ook maar een klein deeltje van uit maakt. Ik weet dat het niet altijd gemakkelijk geweest is om voor mij te zorgen, maar jullie hebben nooit opgegeven, en het nestje dat jullie daarmee bouwden zal voor altijd nauw aan het hart liggen. Grote broer, Bert, ik ben blij dat wij zo'n dichte band hebben, en de tijd wijst alleen maar in een richting die deze band ziet versterken. Een speciale dank om doorheen dit doctoraatsverhaal met mij samen te leven; ik beloof dat ik het huishoudelijk onevenwicht — vooral tijdens die laatste loodjes, zal vereffenen ;).

Wouter Vervust

September 2024

# TABLE OF CONTENTS

# SUMMARY

Molecular dynamics (MD) simulations are a powerful tool to study biological processes at the molecular level. While modern high performance computational infrastructures allow simulations to probe millisecond long simulations, many biological processes extend largely beyond this timescale. In this thesis, techniques are developed for kinetic analysis of long-timescale biological processes using MD simulations.

**Chapter 1** provides an introduction to the this work. It highlights the growth of MD simulations towards larger temporal and spatial scales, emphasizing that increasing the temporal scale is an especially daunting task. The problem is magnified when kinetic analysis is of interest, where the observation of numerous transitions is required for reliable statistics. The chapter discusses how large energy barriers can result in ergodic sampling problems, where long waiting times in metastable states hinder the sampling of other relevant configurations. While advancements in facilitated phase space exploration and subsequent free energy reconstruction techniques have been prolific, these methods typically do not allow accurate analysis of kinetics. To address this shortcoming, path sampling methodologies have been developed. These methods focus on statistical ensembles of trajectories rather than configurations, enabling the study of rare and/or slow events and the calculation of their rate constants. In the final section, the specific research goals of the dissertation are outlined. The work is divided into two main parts. The first part is dedicated to improvements on path sampling methodologies (Chapters 2-6), while the second part develops a diffusive network model to specifically investigate slow oxygen kinetics within myelin sheaths (Chapter 7). The chapter concludes with an outline of the dissertation, detailing the contributions of each subsequent chapter towards achieving the research objectives.

**Chapter 2** provides the theoretical background of path sampling methodologies. It introduces the concepts of path space and path en-

sembles, with a specific focus on Transition Interface Sampling (TIS) and Partial Path Transition Interface Sampling (PPTIS). While TIS provides exact rate constants by considering paths that retain full memory of the transition, its applicability is limited when metastable states are present along the reaction pathway. PPTIS is introduced to tackle this issue, where a diffusive assumption allows path memory to be truncated by defining path ensembles confined to a local region of phase space. The chapter explains how rate constants can be derived from TIS and PPTIS ensembles, after which the core path sampling algorithm (the shooting move) is explained. The importance of detailed balance in path sampling to ensure proper sampling of the ensembles is then discussed. The chapter concludes by discussing the incorporation of a replica exchange move in TIS, which significantly improved ergodic sampling, and how an implementation of this formalism is lacking for PPTIS.

**Chapter 3** delves into the Replica Exchange Partial Path Transition Interface Sampling (REPPTIS) methodology, enhancing the traditional PPTIS framework by incorporation of a replica exchange move. This addition improves ergodic sampling, where the extension and subsequent exchange of paths between ensembles facilitates the exploration of otherwise inaccessible regions of phase space. The effectiveness of REPPTIS is demonstrated through the study of two model systems for membrane permeation, where the resulting rate estimations display improved accuracy compared to PPTIS. The chapter provides the algorithmic implementation of the replica exchange move in the PyRETIS software package, which allows for coupling with various MD simulation engines.

**Chapter 4** introduces an analysis framework to estimate full path lengths and rate constants from PPTIS and REPPTIS simulations. It addresses the limitations of (RE)PPTIS in extracting time-dependent properties such as fluxes, rates, and mean first passage times. By considering a long MD trajectory as a chain of overlapping PPTIS path segments, a Markov state model (MSM) is constructed where the states represent specific path types of the PPTIS ensembles. The transition matrix of this model is given directly by the local crossing probabilities of the PPTIS ensembles. As such, an estimation of full path lengths can be made by reconstructing them as a random walk in the MSM. Additionally, REPPTIS is applied to study the dissociation kinetics of the benzamidine molecule from the trypsin protein. Sampling of the ensembles close to the bound state is done inefficiently, where the replica exchange move is seen to be infrequently

performed. Using a long MD simulation to probe the bound state, the rate estimate is found to be in close agreement with experimental values.

**Chapter 5** examines the dissociation kinetics of imatinib from the ABL kinase domain and its mutated variants. Initially, the RETIS method was used, but it proved infeasible due to the presence of long-lived metastable states along the dissociation pathways, leading to an extremely low acceptance rate. Subsequent REPPTIS simulations faced significant challenges in achieving convergence for reliable rate estimates due to a dense network of metastable states that are separated by high energy barriers orthogonal to the order parameter. This introduced hidden timescale separations that the methodology could not overcome. These issues highlight the limitations of using a one-dimensional order parameter to capture the complex energy landscape of the ABL-imatinib system. The study concludes with the need for multi-dimensional order parameters to better handle metastable states and their associated energy barriers, underscoring the complexity of accurately modeling drug dissociation kinetics in biological systems.

**Chapter 6** introduces a novel extended REPPTIS (REPPEXTIS) path sampling methodology aiming to further improve path memory and ergodic sampling. REPPEXTIS paths are based on the MSM framework introduced in Chapter 4, where its paths are chains of continuously connected PPTIS path segments. The chapter begins by defining the new path ensembles, after which the shoot-and-extend move to sample its paths is explained, and a proof for detailed balance is given. REPPEXTIS paths are contained within multiple ensembles, and can therefore be swapped without the need for additional MD integration. As such, the infinite swapping formalism of $\infty$RETIS is adapted and incorporated in the methodology. It is then described how the use of multiple replicas could enhance replica exchange even further. The methodology is then applied to a shuffleboard potential to demonstrate promising results in path memory increase. An application to a two-dimensional rugged free energy surface is then performed to test the multiple-replica formalism. While early results are promising, there remains much to be explored in future work. The chapter concludes by discussing a proposal for a multi-dimensional REPPEXTIS implementation.

**Chapter 7** investigates how myelination affects oxygen storage and transport in myelin. A diffusive network model is constructed to

study the slow oxygen kinetics within myelin sheaths. MD simulation data of oxygen in one phospholipid bilayer is artificially enlarged to a represent myelin as a stack of multiple bilayers. It is demonstrated that stacked phospholipid membranes increase oxygen storage capacity, suggesting that myelin may act as an oxygen reservoir for nearby oxygen-consuming mitochondria in axons. This enhancement of oxygen storage is seen to level off for large amounts of bilayers, and is diminished by the presence of additional water layers as seen in some cancer cells. The chapter than discusses how the model was extended with time-dependent boundary conditions, and how (un)loading of the stored oxygen within a bilayer follows first-order kinetics. By defining membrane resistance (to oxygen permeation) and membrane capacity (for oxygen storage) parameters, an electric RC circuit analogy is constructed that intuitively captures the time-dependent storage and release of oxygen in one bilayer. The chapter then explains how this model can be expanded into a ladder circuit of RC elements to model oxygen transport through myelin sheaths. The results show that storage capacity and resistance increase linearly with the number of bilayers in the stack, whereas the response time to changes in oxygen concentration displays a quadratic increase. Embedding the ladder circuit within extracellular and axonal solvent compartments, the model can be used to investigates oxygen transport from capillaries to axonal mitochondria. This is done to investigate the effect of myelination on oxygen transport during neuronal activity, which is characterized by a sudden increase of oxygen consumption rate. It is then demonstrated that higher levels of myelination can extend the duration of sustained increased neuronal activity, where the strength of this enhancement depends largely on the onset time of oxygen consumption rate increase and the resting axonal oxygen concentration. While a multitude of scenarios is considered, none of the configurations can sustain increased oxygen demand for much longer than one hundred microseconds, underscoring the functional nature of the hyperemic cerebral blood flow response to restore oxygen levels.

**Chapter 8** concludes the dissertation by summarizing the main findings and contributions of the work, and discussing potential future research directions.

# SAMENVATTING

Moleculaire dynamica (MD) simulaties zijn een krachtig hulpmiddel om biologische processen op moleculair niveau te bestuderen. Terwijl moderne computationele rekenkracht simulatietijden van milliseconden lang toelaten, reiken veel biologische processen ver voorbij deze tijdschaal. In dit proefschrift worden technieken ontwikkeld voor kinetische analyse van biologische processen op lange tijdschaal met behulp van MD-simulaties.

**Hoofdstuk 1** biedt een introductie tot dit werk. Het benadrukt de groei van MD-simulaties naar grotere tijd- en ruimtelijke schaal, met de nadruk op het feit dat het vergroten van de temporele schaal een bijzonder moeilijke taak is. Het probleem wordt vergroot wanneer kinetische analyse van belang is, waarbij de observatie van talrijke transities nodig is om betrouwbare statistiche uitspraken te doen. Het hoofdstuk bespreekt hoe grote energiebarrières kunnen resulteren in ergodische bemonstering problemen, waarbij lange wachttijden in metastabiele toestanden het bemonsteren van andere relevante configuraties belemmeren. Hoewel er veel vooruitgang is geboekt in het faciliteren van fase-ruimte verkenning en de daarbijhorende technieken voor de reconstructie van vrije energie, staan deze methoden doorgaans geen nauwkeurige analyse van kinetiek toe. Om deze tekortkoming aan te pakken, zijn padbemonsteringmethoden ontwikkeld. Deze methoden richten zich op statistische ensembles van trajecten in plaats van configuraties, waardoor de studie van zeldzame en/of langzame gebeurtenissen en de berekening van hun snelheidsconstanten mogelijk wordt. In het laatste deel van dit hoofdstuck worden de specifieke onderzoeksdoelen van het proefschrift geschetst. Het werk is verdeeld in twee hoofdonderdelen. Het eerste deel is gewijd aan verbeteringen in padbemonsteringmethoden (Hoofdstukken 2-6), terwijl het tweede deel een diffuus netwerkmodel ontwikkelt om specifiek langzame zuurstofkinetiek binnen myelinescheden te onderzoeken (Hoofdstuk 7). Het hoofdstuk eindigt met een overzicht van het proefschrift, waarin de bijdragen

van elk volgend hoofdstuk aan het bereiken van de onderzoeksdoelen worden beschreven.

**Hoofdstuk 2** biedt de theoretische achtergrond van padbemonsteringmethoden. Het introduceert de concepten van padruimte en padensembles, met specifieke nadruk op Transition Interface Sampling (TIS) en Partial Path Transition Interface Sampling (PPTIS). Terwijl TIS exacte snelheidsconstanten biedt door paden te beschouwen die hun volledige geschiedenis van de overgang behouden, is de toepasbaarheid ervan beperkt wanneer metastabiele toestanden aanwezig zijn langs het reactiepad. PPTIS wordt geïntroduceerd om dit probleem aan te pakken, waarbij een diffuse veronderstelling het toestaat om de padgeschiedenis in te korten door padensembles te definiëren die beperkt zijn tot een lokaal gebied van de fase-ruimte. Het hoofdstuk legt uit hoe snelheidsconstanten kunnen worden afgeleid van TIS- en PPTIS-ensembles, waarna het kernpadbemonsteringsalgoritme (de 'shooting move') wordt uitgelegd. Het belang van gedetailleerde balans in padbemonstering om een juiste bemonstering van de ensembles te waarborgen wordt vervolgens besproken. Het hoofdstuk eindigt met een bespreking van de positieve impact van replica-uitwisseling in TIS, die de ergodische bemonstering aanzienlijk verbetert, en hoe een implementatie van zo een uitwisselingsformalisme ontbreekt voor PPTIS.

**Hoofdstuk 3** gaat dieper in op de Replica Exchange Partial Path Transition Interface Sampling (REPPTIS) methodologie, die het traditionele PPTIS-framework verbetert door de invoering van een replica-uitwisselingsbeweging. Deze toevoeging verbetert de ergodische bemonstering, waarbij de verlenging en daaropvolgende uitwisseling van paden tussen ensembles de verkenning van anders ontoegankelijke gebieden van fase-ruimte vergemakkelijkt. De effectiviteit van REPPTIS wordt aangetoond door de studie van twee modelsystemen voor membraanpermeatie, waarbij de resulterende snelheidsschattingen verbeterde nauwkeurigheid vertonen in vergelijking met PPTIS. Het hoofdstuk biedt de algoritmische implementatie van de replica-uitwisselingsbeweging in het PyRETIS-softwarepakket, dat koppeling met verschillende MD-simulatiepakketten mogelijk maakt.

**Hoofdstuk 4** introduceert een analysekader om volledige padlengtes en snelheidsconstanten te schatten uit PPTIS- en REPPTIS-simulaties. Het behandelt de beperkingen van (RE)PPTIS voor het bepalen van tijdsafhankelijke eigenschappen zoals fluxen, snelheden en gemiddelde eerste passage tijden. Door een lange

MD-traject als een keten van overlappende PPTIS-padsegmenten te beschouwen, wordt een Markov-staatmodel (MSM) geconstrueerd waarbij de staten specifieke padtypen van de PPTIS-ensembles vertegenwoordigen. De overgangsmatrix van dit model wordt direct gegeven door de lokale doorkruisingsprobabiliteiten van de PPTIS-ensembles. Zo kan een schatting van volledige padlengtes worden gemaakt door ze te reconstrueren als een willekeurige wandeling in de MSM. Daarnaast wordt REPPTIS toegepast om de dissociatiekinetiek van het benzamidinemolecuul uit het trypsine-eiwit te bestuderen. De bemonstering van de ensembles dichtbij de gebonden toestand wordt inefficiënt uitgevoerd, waarbij de replica-uitwisselingsbeweging zelden wordt uitgevoerd. Met behulp van een lange MD-simulatie om de gebonden toestand te onderzoeken, wordt een snelheidschatting gevonden die nauw overeenkomt met experimentele waarden.

**Hoofdstuk 5** onderzoekt de dissociatiekinetiek van imatinib uit het ABL-kinasedomein en zijn gemuteerde varianten. Aanvankelijk werd de RETIS-methode gebruikt, maar deze bleek onuitvoerbaar vanwege de aanwezigheid van langdurige metastabiele toestanden langs de dissociatiepaden. Latere REPPTIS-simulaties ondervonden aanzienlijke uitdagingen om convergentie voor betrouwbare snelheidschattingen te bereiken vanwege een dicht netwerk van metastabiele toestanden die worden gescheiden door hoge energiedrempels loodrecht tot de ordeparameter. Dit introduceerde verborgen tijdschaal scheidingen die de methodologie niet kon overwinnen. Deze problemen benadrukken de beperkingen van het gebruik van een eendimensionale ordeparameter om het complexe energielandschap van het ABL-imatinib systeem vast te leggen. De studie concludeert met de noodzaak van multidimensionale ordeparameters om metastabiele toestanden en hun bijbehorende energiedrempels beter aan te kunnen, wat de complexiteit benadrukt van het nauwkeurig modelleren van medicijndissociatiekinetiek in biologische systemen.

**Hoofdstuk 6** introduceert een nieuwe uitgebreide REPPTIS (REPPEXTIS) padbemonsteringmethodologie ontworpen om padherinnering en ergodische bemonstering te verbeteren. REPPEXTIS-paden zijn gebaseerd op het MSM-formalisme geïntroduceerd in Hoofdstuk 4, waarbij de paden ketens zijn van continu geëxtendeerde PPTIS-padsegmenten. Het hoofdstuk begint met het definiëren van de nieuwe padensembles, waarna de schiet-en-verlengbeweging om de paden te bemonsteren wordt uitgelegd, en een bewijs voor

gedetailleerde balans wordt gegeven. REPPEXTIS-paden zijn vervat in meerdere ensembles, en kunnen daarom worden uitgewisseld zonder de noodzaak van additionele MD-integratie. Zo wordt de oneindige uitwisselingsformalisering van $\infty$RETIS aangepast en opgenomen in de methodologie. Vervolgens wordt beschreven hoe het gebruik van meerdere replica's de replica-uitwisseling nog verder zou kunnen verbeteren. De methodologie wordt vervolgens toegepast op een shuffleboard-potentiaal om de toename van padherinnering aan te tonen, en op een tweedimensionaal ruw vrije-energiesoppervlak om de multiple-replica-formalisering te testen. Hoewel vroege resultaten veelbelovend zijn, blijft er nog veel te verkennen in de toekomst. Het hoofdstuk eindigt met de bespreking van een voorstel voor een multidimensionale REPPEXTIS-implementatie.

**Hoofdstuk 7** onderzoekt hoe myelinisatie zuurstofopslag en -transport in myeline beïnvloedt. Een diffuus netwerkmodel wordt geconstrueerd om de langzame zuurstofkinetiek binnen myelinescheden te bestuderen. MD-simulatiegegevens van zuurstof in één fosfolipide dubbellaag worden kunstmatig vergroot om myeline te vertegenwoordigen als een stapel van meerdere dubbellagen. Het wordt aangetoond dat gestapelde fosfolipide membranen de zuurstofopslagcapaciteit verhogen, wat suggereert dat myeline kan fungeren als een zuurstofreservoir voor nabijgelegen zuurstofverbruikende mitochondriën in axonen. Deze verbetering van zuurstofopslag blijkt af te vlakken voor grote hoeveelheden dubbellagen en wordt verminderd door de aanwezigheid van extra waterlagen zoals te zien in sommige kankercellen. Het hoofdstuk bespreekt vervolgens hoe het model is uitgebreid met tijdsafhankelijke randvoorwaarden, en hoe het (ont)laden van de opgeslagen zuurstof binnen een dubbellaag een eerste-orde kinetiek proces is. Door membraanweerstand (tegen zuurstofpermeatie) en membraancapaciteit (voor zuurstofopslag) parameters te definiëren, wordt een elektrische RC-circuitanalogie geconstrueerd die op intuïtieve wijze de tijdsafhankelijke opslag en afgifte van zuurstof in één dubbellaag weergeeft. Het hoofdstuk legt vervolgens uit hoe dit model kan worden uitgebreid tot een laddercircuit van RC-elementen om zuurstoftransport door myelinescheden te modelleren. De resultaten tonen aan dat opslagcapaciteit en weerstand lineair toenemen met het aantal dubbellagen in de stapel, terwijl de responstijd op veranderingen in zuurstofconcentratie een kwadratische toename vertoont. Door het laddercircuit in te bedden binnen extracellulaire en axonale solvent fasen, kan het model worden gebruikt om zuurstoftransport van haarvaten naar axonale mitochondriën te onderzoeken. Dit wordt gedaan om het effect van myelinisatie op zuurstoftransport tijdens neuronale

activiteit te onderzoeken, die wordt gekenmerkt door een plotselinge toename van de zuurstofverbruikssnelheid. Vervolgens wordt aangetoond dat hogere niveaus van myelinisatie de duur van aanhoudende verhoogde neuronale activiteit kunnen verlengen, waarbij de sterkte van deze verbetering grotendeels afhangt van de begintijd van de toename van de zuurstofverbruikssnelheid en de equilibrium axonale zuurstofconcentratie. Hoewel een groot aantal scenario's wordt overwogen, kan geen van de configuraties verhoogde zuurstofvraag veel langer dan honderd microseconden ondersteunen, wat de functionele aard van de hyperemische cerebrale bloedstroomrespons benadrukt om zuurstofniveaus te herstellen.

**Hoofdstuk 8** sluit het proefschrift af door de belangrijkste bevindingen en bijdragen van het werk samen te vatten en mogelijke toekomstige onderzoeksrichtingen te bespreken.

# I

## Part I - Advanced Molecular Dynamics Simulation Techniques for Kinetic Analysis of Biological Systems

# INTRODUCTION

The Digesting Duck (Fig. 1.1) was a mechanical automaton built by Jacques de Vaucanson, exhibited to the public in the winter of 1738 in Paris [1]. An intricate piece of machinery, reminiscent of a clockwork, allowed the duck to make a series of lifelike movements, quack, gurgle water, and even eat grain from your hand, 'digest' it, and after a pause, 'excrete' it [2]. De Vaucanson's automata received much appraisal, in which he was hailed 'Prometheus's rival' by Voltaire, who further deemed the Digesting Duck 'the sole reminder of France's



**Figure 1.1:** Left: the Digesting Duck, an automaton built by Jacques de Vaucanson in the eighteenth century. Right: ball-and-stick representation of duck $\delta 1$ crystallin, an eye lens protein. (PDB entry 1HY0 [3, 4]).

glory' [1]. Its success can be attributed to the dominant mechanically reductionistic view of the world at the time, finding its roots in Descartes' highly influential *Meditations*, published approximately a century prior to the clockwork duck (1641) [5, 6]. Not much later, Newton published his *Principia Mathematica* in 1687. The ordered, rational, and mathematical description of celestial body motion was the culmination of the Scientific Revolution, where once the heavens were considered the realm of the divine; now they were governed by the same laws as the Earth. The *empirical method* of Bacon (which Newton was inspired by [7]) had gained traction, where systematic observation and experimentation had now become the cornerstones of scientific discovery. It is in this clockwork universe of Newton and Descartes (mechanical laws for terrestrial as well as celestial objects) that scientists, armed with the newly developed scientific method (observation and experimentation), developed reductionistic theories, attempting to 'reduce' complex systems into the Newtonian paradigm [8].

Hence, ducks became Digesting Ducks, because when it quacks like a duck, walks and talks like a duck, it *is* a duck? If you failed to answer this rhetorical question, you can find inspiration in Aristotle's writings 2000 years prior to our Digesting Duck. In his *Metaphysics*, pondering over "the cause of a substance's unity", Aristotle notes that "totality is not, as it were, a mere heap, but the whole is something beside the parts" [9]. The more widespread variant states that "the whole (system) is more than just the sum of its parts", acknowledging the importance of 'emergent phenomena' by the complex interplay of a system's constituents. An example can be found in water ($H_2O$), whose molecular structure cannot be predicted by the atomic properties of hydrogen and oxygen (a molecular theory is required to explain bonding). Macroscopic properties of water, like surface tension, cannot be predicted by the molecular structure alone (a statistical mechanical theory is required to explain the collective behavior of $H_2O$ molecules).

However, there need not be a 'great divide' between reductionist and holistic science (i.e. anti-reductionistic science, often wrongly associated with pseudo-science), as both approaches are interdependent and complementary [10]. Nevertheless, reductionism has been (and continues to be) highly successful, where a prime example is found in the molecularization of biology. Molecular biology has not only revolutionized biology, but also the medicinal field, where knowledge of genetic code (the small) can sometimes be used to both diagnose

and treat diseases (the large). As Richard Dawkins once said: "Reductionism is one of those things, like sin, that is only mentioned by people who are against it" [11].

Science has tremendously progressed since the clockwork universe of Newton and Descartes, where a cog-and-gear representation for biological life is now mostly found in museums, the 'steampunk' section in bookstores, and introductory texts of PhD dissertations. A more advanced clockwork that fits the Digital Age will be used in this dissertation: molecular dynamics (MD) simulations. MD simulations provide a computational lens of atomic resolution, visualizing biological systems by simulating the motion of its constituent atoms. While simulating an entire duck by its atomic content is not feasible today or in the nearby future, small parts of the duck can be tackled, where a typical MD simulation could lay bare the motions and interactions of duck proteins (Fig. 1.1B). While the MD approach is obviously reductionistic (biological system = a collection of atoms), MD simulations provide a holistic view by letting these atoms interact, from which emergent properties can be observed. An example is found in the recent (2019) MD simulation of a 136 million atom-scale model of an entire cell organelle, which was able to reproduce phenotypic properties from atomistic detail [12]. The size of this organelle fluctuated around 50 to 60 nm, and the MD simulated time was 0.5 μs, which remain very impressive numbers even 6 years later.

This work addresses the challenge of achieving biologically relevant sizes and timescales in molecular dynamics (MD) simulations, crucial for understanding the kinetics of biological processes. It makes significant contributions in two areas. Firstly, it enhances methodologies for extracting long-timescale kinetics through both algorithmic improvements and the introduction of novel approaches. Secondly, a method that artificially upscales bilayer permeation simulations is introduced, which is applied to investigate the previously unexplored kinetics of oxygen in myelin sheaths.

Before delving into these contributions, the specific nature of these challenges is outlined in the following introductory sections. First, the field of biomolecular modeling is introduced in Sec 1.1, highlighting that despite significant advances in computational power, many biological processes remain inaccessible. The characteristics of biological free energy landscapes are discussed in Sec 1.2, along with how they can be explored by enhanced sampling techniques. It is then demonstrated that these techniques are generally unsuitable for kinetic analysis, prompting the development of path sampling methods,

discussed in Sec 1.3. The structure of this work is presented in Sec 1.4, together with the research objectives. Finally, Sec. 1.5 lists the papers that constitute this work, where three are published (**papers I, II, IV**) and two are in preparation (**papers III, V**).

## 1.1 MOLECULAR SIMULATION OF BIOLOGICAL SYSTEMS

Molecular simulation can concisely be defined as a computational manifestation of statistical mechanics [13], which provides the theoretical framework to derive macroscopic properties from a microscopic description of a many-body system. Bridging microscopic and macroscopic worlds is crucial, especially given the constraints of experimental techniques in probing atomic-scale dynamics [14]. Cornerstone methods of molecular simulation include molecular dynamics (MD) and Monte Carlo (MC) simulations, whose core algorithms have remained fairly constant since their inception in the 1950s [14–16]. Both methods result in an ensemble of configurations, where ensemble averages of observables can be calculated and compared to experimental data. In MC simulations, configurations are generated randomly according to the underlying probability distribution, whereas MD simulations generate configurations by integrating Newton's equations of motion. The output of an MD simulation is essentially a movie, as it captures the trajectories of all the atoms, effectively serving as a computational microscope. However, the advantageously high resolution of MD simulations naturally introduces challenges for their applicability to the large and slow world of biology.

These terms, 'large' and 'slow', are best understood through examples, starting with the spatial scale. The first biological MD simulations in the 1970s dealt with small proteins consisting few hundred atoms [17, 18]. Computational power has increased exponentially for more than five decades [19, 20], which has allowed the MD system size to grow to an entire satellite tobacco mosaic virus ($\sim$1 million atoms) in 2006 [21], to an entire gene locus ($>$1 billion atoms) in 2019 [22] and an entire SARS-CoV-2 virus ($\sim$305 million atoms) in 2021 [23]. While breaking the barrier of 1 billion atom simulations is a tremendous achievement, it is still a far cry from what many would consider the basic unit of life: the cell [20]. A mammalian cell with a typical cell volume of $2000\,\mu m^3$- $4000\,\mu m^3$ contains approximately $10^{10}$ proteins [24], equivalent to $4 \times 10^{13}$ atoms when taking the average protein to contain 4000 atoms. This only comprises the protein content ($\sim 20\%-30\%$ volume fraction [25]), and inclusion of the solvent would

result in approximately $10^{14}$ (100 trillion) atoms per cell. This is 5 orders of magnitude beyond the largest MD simulations ever performed. Predicting when this will be possible is difficult, as (1) Moore's law is reaching its end [20], (2) artificial intelligence (AI) enhanced simulations will become stronger [26], and (3) whether or not the community deems it a goal worth pursuing [27]. A step towards larger systems is provided by coarse-grained (CG) models, where clusters of atoms are reduced to a single 'bead', significantly reducing the degrees of freedom (DoFs) [28, 29]. This loss of atomic precision results in system size gain and longer simulations, which recently allowed the simulation of an entire *minimalistically engineered* cell ($\sim 550$ million coarse-grained beads $\approx 6$ billion atoms $\ll$ mammalian cell) [30]. Periodicity in biological systems can be exploited to artificially enlarge them, which is used in **papers IV and V** to examine oxygen kinetics in myelin sheaths. By viewing myelin as a stack of phospholipid bilayers, data from a single-bilayer MD simulation proved adequate for the study. The study further discusses how this methodology can be extended to investigate the permeation of other small molecules, where the calculation of a few single-bilayer properties can be used to predict kinetics in a myelin sheath.

For the temporal scale, simulation times grew from 9.2 ps for the bovine pancreatic trypsin inhibitor (458 'pseudo' atoms) in 1977 [17, 18], to breaking the microsecond barrier in 1998 for the HP-36 villin headpiece subdomain (36 residues and $\sim$3000 water molecules) [31], to breaking the millisecond barrier in 2009 for the BPTI protein (17,758 atoms) [32]. This last example ran on Anton, a special-purpose supercomputer with hardware specifically designed for MD simulations [33]. A second and third version of Anton have since been introduced, where Anton3 boasts 100 µs per day for a million atom system [34, 35]. While access to the millisecond timescale is highly impressive, this is dwarfed by many biological processes of interest. Consider, for example, the dissociation of the drug molecule imatinib from the kinase protein ABL, with a dissociation rate of approximately $0.001 \, \text{s}^{-1}$ [36]. One obtains a simulation rate of $\sim$300 ns/day after solvating the protein-drug complex into a simulation box of 50 000 atoms, on a recent A100 GPU with an AMD EPYC CPU on the Gromacs MD simulation software (version 2021.3) [37]. Assuming dissociation is a Poisson process, an unbinding event is expected to occur after 9126 millennia of simulation time. Even on the priceless Anton3 supercomuter, one would require a *few* millennia. Bridging this challenging gap between the timescale of biological processes and

the time accessible to MD simulations requires the development of advanced methodologies, which is the topic of **papers I-III**.

It must be noted that, to create a *continuous* (or *conventional*) trajectory, increasing the system size (atom count $N$) can be considered 'easier' than increasing the timescale (simulation time $T$). This can be surprising, noting that computational complexity scales as $\mathcal{O}(N \log N \times T)$ or $\mathcal{O}(N \times T)$, depending on usage of particle mesh Ewald or fast multipole methods for the long-range interactions, respectively [14, 38–42]. Larger and larger simulation boxes become possible as spatial parallelization schemes (e.g. domain decomposition [43]) can efficiently distribute (spatial) computational load over High Performance Compute (HPC) platforms with low communication latency, where individual compute nodes deal with only a fraction of the total atoms $N$ [44]. While approximative methods exist [45–47] there is no such thing as 'temporal domain decomposition' for conventional MD simulations, as the chaotic nature of many-body systems enforces sequential integration. In short, the spatial size can be increased by using *more* hardware, whereas the temporal reach of MD simulations requires *faster* hardware. Algorithmic advancements have also improved conventional simulation rates, mainly by increasing the integration time step which is normally limited to a $\leq$ femtosecond by the high-frequency motion of covalently bonded hydrogens (NH, CH, OH, etc.) [48, 49]. This can be increased to 2 fs by keeping H-bonds rigid [50–52], upwards to 4 fs by mass repartitioning schemes [53], and upwards to 5 fs by treating hydrogens as massless virtual interaction sites [54, 55]. All of the algorithmic advances and (the ability to use) hardware advancements are readily available to the scientific community in (bio)molecular software packages like Gromacs [37], NAMD [56], AMBER [57], CHARMM [58], ACEMD [59], and OpenMM [60]. While of paramount importance to the field, the development of increasingly accurate force fields is not the focus of this work, and the reader is referred to Refs. [61–64] for recent reviews.

The attainable timescales of MD simulations are compared with the timescales of relevant biological processes in Fig. 1.2. It also shows the accessible timescales for common experimental methods to probe protein motions.

**Figure 1.2:** Timescales of protein and phospholipid bilayer motions span many orders of magnitude. Timescales attainable by experimental techniques are also shown. H/D exchange: hydrogen-deuterium exchange, AFM: atomic force microscopy, FRET: fluorescence resonance energy transfer, NMR: nuclear magnetic resonance, IR: infrared spectroscopy, CHOL: cholesterol. Values are collected from Refs. [65–75].

Two observations can be made from Fig. 1.2. Firstly, the timescales of biological processes span many orders of magnitude, where important processes such as slow conformational changes, protein (un)folding, protein-drug (un)binding and membrane-drug permeation can happen at timescales well out of reach for MD simulations. Secondly, a distinction in MD accessible timescales is made between 'sampling' and 'kinetics'. The sampling timescale denotes the accessible length of MD trajectories and reaches upwards to milliseconds. As kinetic analysis of a process requires the observation of *hundreds* of events, the accessible timescale for kinetic assessment via MD simulations is a few orders of magnitude shorter than the sampling timescale.

Clearly, there is a need for methods that push our computational microscope to probe kinetics at longer timescales. A first step towards kinetic analysis is the use of enhanced sampling techniques

that facilitate phase space exploration and reconstruction of free energy profiles.

## 1.2 FREE ENERGY BARRIERS AND ENHANCED SAMPLING

Fig. 1.3 shows the free energy of a system as a function of a reaction coordinate (RC) $\lambda$, i.e. a Landau-type free energy [76]. The RC is typically chosen to be representative of the process of interest, such as the distance between a drug molecule and a protein binding pocket. From an MD simulation, a property of interest $A$ can be



**Figure 1.3:** Free energy profile along a reaction coordinate $\lambda$. Thermodynamics involves the energy differences between states (e.g. $\Delta F$), while kinetics involves the energy barriers separating states (e.g. $\Delta^{\mp}F$ and $\Delta^{\pm}F$). A change of the system, such as a mutated residue in a protein, can severely alter the thermodynamics and kinetics of the system (black versus red curves). This is discussed in **Chapter IV** for the dissociation of the drug molecule imatinib from the ABL kinase protein and mutated variants thereof.

calculated by averaging over the simulation trajectory, resulting in a *time averaged* value $\bar{A}$. The ergodic hypothesis states that, for an infinitely long trajectory, the time averaged value equals the *ensemble averaged* value $\langle A \rangle$, where the property is averaged over all possible configurations of the system. While this is a powerful concept, it is not always true. Even worse, MD trajectories are not infinitely long, where the time averaged value can become highly dependent on the initialization of the system. For a temperature controlled ergodic MD simulation, all configurations should be sampled according to the Boltzmann distribution, allowing for a reconstruction of the entire free energy profile. To overcome energy barriers, all the DoFs of the system must align favorably to provide the necessary energy to *activate* the process, where the activation energy dependency was put into equation as early as 1889 by Arrhenius [77]. Large energy barriers, or a large collection of smaller barriers, can severely hinder *ergodic*

*sampling.* For example, an MD simulation initiated on the left of the large energy barrier in Fig. 1.3 will likely not sample the (dominant!) right energy well within accessible simulation time. Throughout this work, ergodicity will be used to denote the ability of a simulation to sample all relevant states of a system within reasonable simulation time (and not the strict mathematical definition). This problem of non-ergodic sampling seeded the development of enhanced sampling strategies that aim to facilitate the exploration and subsequent reconstruction of phase space.

A first class of methods changes the Hamiltonian along a (low-dimensional) set of predefined collective variables (CVs). The oldest (but still very popular) method is umbrella sampling [78]. Newer methods use a memory dependent potential to disfavor previously explored regions of phase space, such as conformational flooding [79], local elevation, [80], metadynamics [81, 82], adaptive biasing force [83, 84], and variationally enhanced sampling [85]. Analysis based on the weighted histogram analysis method (WHAM [86]) or the multistate Bennett Acceptance Ratio (MBAR [87]) are subsequently used to recover the true free energy profile. Often the reaction mechanism is not well-known, and important energy barriers remain orthogonal to the predefined CVs, hindering ergodic sampling. A biasing strategy that does not require a predefined CV is accelerated MD [88], where biasing is based on a potential energy threshold.

A second class of methods achieves enhancement not by adding a biasing force, but by altering the underlying canonical probability distribution itself. These are replica exchange methods (parallel tempering, also without CV) [89–91] and variants thereof, such as simulated tempering [92], Wang-Landau sampling [93] and integrated tempering sampling [94].

Once the free energy profile is known, Transition State Theory (TST) can be used to calculate reaction rate constants [95–97]. TST, the theoretical framework of activated processes that improved on Arrhenius's work, makes strong assumptions about the reaction mechanism (no recrossings, equilibrium of reactant and transition state, RC with lowest saddle point), where its rate estimations can be highly inaccurate. To this end, the reactive flux (RF) formalism (Bennett-Chandler approach [98]) was introduced to provide a dynamical correction factor accounting for TST inaccuracies. This methodology essentially probes the dynamics of crossing the TST, by launching short MD trajectories from a dividing surface and tracking whether they commit to the reactant or product.

The RF formalism produced impressively accurate rate constants for low-dimensional systems such as small chemical reactions, where the transition state (ensemble) can be approximated by a few saddle points on a smooth free energy landscape. For biological systems, however, it proved to be hardly applicable. The root cause lies in their 'rough' or 'rugged' free energy landscapes, containing numerous local minima separated by high barriers (Fig. 1.4). This shape is typical for condensed matter systems, where the activated process is determined by the favorable alignment of hundreds of weak molecular interactions [76, 99].



**Figure 1.4:** A 2D representation of a 'rugged' or 'rough' free energy landscape often encountered in biological systems.

Good RCs and dividing surfaces on such rugged free energy surfaces are elusive, where bad choices results in very poor rate estimates (by TST) that are difficult to correct (with RF). This is especially the case for slow reactions, where the 'short' RF trajectories become infeasibly long.

## 1.3  PATH SAMPLING AND KINETICS

The shortcomings of TST and RF led to the development of strategies that directly tackle the kinetics of the system, rather than trying to derive them from the static (thermodynamics) free energy profile. This required a shift from statistical ensembles of configurations to a theory of statistical ensembles of trajectories. Due to their importance to the methodologies developed in **papers I-III**, the theory of path ensembles is discussed in more detail in chapter 2.

The first path sampling method was Transition Path Sampling (TPS), which generates reactive trajectories without the need for an

RC [100–106]. It achieves this in MC fashion, where new MD traject-
ories are generated from old ones using the so-called *shooting* move.
This involves perturbing the momenta of a randomly chosen phase
point of an old trajectory, after which it is propagated forwards and
backwards in time in the hope of generating a new reactive trajectory.
Metropolis acceptance rules for these paths enforce detailed balance
such that the correct path ensemble is sampled. The classical TPS
algorithm saw broad application, including studies of conformational
changes of proteins [107–109] and nucleotides [110, 111].

In the $\sim$25 years following the introduction of TPS, many path
sampling methodologies have been developed, where most of them
reintroduced the need for a reaction coordinate or order parameter $\lambda$.

A large improvement was offered by transition interface sampling
(TIS) [112], where partitioning of phase space via interfaces allows a
divide-and-conquer strategy to calculate rate constants as a product
of history-dependent crossing probabilities [113, 114]. A replica ex-
change move was added to TIS that greatly increased the ergodicity
of the method (RETIS [115–117]), which was recently updated with
an infinite swapping protocol ($\infty$RETIS [118, 119]). Multiple state
TIS (MSTIS [120]) was developed to deal with processes other than
the typical reactant-product mechanism.

(RE)TIS is an exact path sampling method, meaning that it gen-
erates paths as if they were generated by a very long MD simulation.
While extremely efficient for rare events, the presence of long-lived
metastable states along the reaction pathways can make application
of (RE)TIS infeasible. For such systems, that are both rare *and* slow,
the Partial Path TIS (PPTIS) method was developed [121]. PPTIS
ensembles are restricted to a local region of phase space, effectively
truncating path memory to allow for shorter paths, where compu-
tational load is reduced at the cost of exactness. While PPTIS can
be applied to a broader range of biophysical systems, the methodo-
logy faces two main challenges. First, the method is more depend-
ent on the choice of $\lambda$ parameter, and second, the path ensembles
do not exchange information, leading to potential ergodicity prob-
lems. These issues are both demonstrated and addressed in **paper I**,
where introduction of a replica exchange move within PPTIS is pro-
posed (REPPTIS [122]). While (RE)PPTIS is a powerful method,
an implementation that couples with common MD engines was never
developed. This is addressed in **paper II**, where the (RE)PPTIS
methodology is implemented in the third version of the PyRETIS
software package [123].

Improvements to TPS and TIS methods have been made by means of new MC moves that aim to increase path acceptance (precision shooting [124], noise-guiding [125], shooting from the top [126] S-shooting [127]), or aim to increase convergence rate and ergodicity by faster path decorrelation ('metadynamics in path space' [128] and the 'web throwing', 'stone skipping' and 'wire fencing' sub-trajectory moves [129, 130]).

Another descendant of TIS is Forward Flux Sampling (FFS [131]), which is based on the non-Metropolis MC scheme of *splitting* to generate new paths [132]. Major advantages of splitting are that (1) it does not require reversible dynamics (as the paths are only propagated forward in time), and (2) it can be used for non-equilibrium systems (as it requires no knowledge of phase point densities) [113]. Disadvantages include larger correlations between generated paths, and a much stronger dependence on the $\lambda$ parameter [133]. Other splitting methods include RESTART [134], weighted ensemble [135] (which was notably used to elucidate the opening dynamics of the SARS CoV-2 spike protein [136]), and adaptive multilevel splitting (AMS) [137].

Another non-Markov chain method is milestoning [138]. Milestoning splits phase space into segments or 'milestones' along a $\lambda$ parameter, from which large amounts of trajectories are spawned. It has become a popular method, especially for reactions covering many metastable states, where combination with Voronoi tessellation allowed (automated) integration of multidimensional $\lambda$ parameters [139–144]. Milestoning and PPTIS, while developed independently, share strong similarities. PPTIS was developed from the TPS viewpoint, to which a $\lambda$ parameter is introduced to segment and probe reactive paths using shorter paths. Milestoning was developed to probe timescales for mechanisms with known $\lambda$ parameters using short trajectories [141]. Milestoning assumes full memory loss when a milestone is hit, while PPTIS uses a softer Markovian approximation where the interface-to-interface transitions are conditionally defined on the previously crossed interface, making it less dependent on the choice of $\lambda$ parameter. On the other hand, milestoning uses time-dependent transition probabilities, while PPTIS does not [113]. Extra memory was introduced in the directional milestoning approach, making it more 'PPTIS-like' [140, 141]. While the (RE)PPTIS method delivered crossing probabilities, it did not yet provide a rate constant and other time-dependent properties. This was addressed in **paper III**, where a Markov State Model (MSM) interpretation of

the PPTIS framework is constructed, making the post-analysis more 'milestoning-like'.

Apart from PPTIS and splitting methods, a third class of methods that uses short trajectories to access longer timescales are MSM based methods [145–149], where macrostates are defined by means of feature selection and clustering algorithms (kinetic coarse-graining [150–153]). Considerable progress has been made via the variational approach to conformation dynamics (VAC [154]) and time-lagged independent component analysis (tlCA or TICA [155]). Further progression was made by bridging variational formulations (allowing the design of algorithms robust in high dimensions) with data-based approaches of machine learning (ML) [156]. A prominent example is VAMPnets [157], presenting a deep-learning approach to molecular kinetics.

## 1.4 Research goals and outline

As outlined in the previous sections, the computational study of many important biomolecular processes requires sampling strategies that go beyond the timescales accessible to conventional MD simulations. Developments in enhanced sampling methods to facilitate phase space exploration and reconstruction of free energy landscapes have been prolific. However, these methods often do not allow an accurate determination of rate constants, where the TST/RF formalism to bridge the static-to-kinetic worlds is hardly applicable to the rough free energy landscapes of biological systems. Assessing kinetics is crucial, however, as the world of biology is a world of dynamics.

Path sampling methodologies bring the powerful mechanism of importance sampling to path space, allowing focused sampling in the rare transition region and subsequent calculation of rate constants. TPS, TIS and RETIS are all in the category of exact path sampling methods, where the resulting trajectories can be seen as if they were extracted from an infinitely long MD simulation. These methods are designed to study *rare events*, which are processes that (1) happen very infrequently on the timescale of molecular fluctuations accessible to MD simulations, and (2) when they do occur, the transition happens rapidly. While (RE)TIS provides an exponential reduction of compute time required compared to MD, it is not without its limitations. If long-lived metastable states are present along the transition pathways, the process is both *rare and slow*, for which (RE)TIS paths become infeasibly long and the method is no longer applicable,

as mentioned in the previous section. PPTIS sacrifices exactness by restricting the path ensembles to a local region of phase space, where the shorter paths reduce computational demand at the cost of inherent approximation by path memory truncation. While this ensures a broader applicability of the PPTIS method, it introduces both ergodicity and accuracy concerns. Improving the PPTIS methodology therefore constitutes the first large research objective of this work.

> **Research Question 1**: *Can the PPTIS methodology be improved to provide reliable kinetic analysis of long-timescale biological processes?*

To achieve this, some intermediate steps are taken. **Chapter II** provides the theoretical background of path sampling, particularly focusing on the path ensembles encountered in TIS and PPTIS simulations. It details how rate constants can be derived from these simulations, and provides a closer look at the core sampling algorithm of these methods: the shooting move.

**Chapters III-VI** contribute directly to Research Question 1 by addressing specific aspects that currently limit the accuracy and applicability of the PPTIS methodology. These limiting aspects are first put into categories, where tackling Research Question 1 can then be broken down into four *intermediate Research Objectives*.

The first category of limitations inherently limits the accuracy of PPTIS, as

**1 a**: PPTIS path ensembles do not exchange information (paths).

- While the replica exchange move was originally proposed for the PPTIS methodology [121], it was only implemented for the TIS method. Ensemble communication via replica exchange can greatly increase ergodic sampling of the path ensembles, where regions that are hardly accessible for one ensemble can be reached via sampling in another ensemble. To address the lack of ensemble communication in PPTIS, the replica exchange is implemented in the PPTIS framework (REPPTIS) in **Chapter III**. To further increase ensemble communication, a novel path sampling method is introduced in **Chapter VI**.

**1 b**: PPTIS path ensembles have limited memory.

- Due to the restriction PPTIS ensembles to a local region of phase space, paths no longer 'remember' whether they are reactive or non-reactive, or whether they come from the reactant rather than product state. This theoretically limits the accuracy of the rate reconstruction, where a good choice of $\lambda$ parameter and subsequent phase space partitioning via interfaces becomes crucial. The novel path sampling methodology addresses this issue in **Chapter VI**, where path memory is increased by extending paths to sample outside the local region.

The second category pertains to the applicability of PPTIS, as

**1 c**: PPTIS lacks a framework to extract time-dependent properties

- As a method to study kinetics, (RE)PPTIS has no theoretical framework to directly calculate fluxes, rates and mean first passage times from the simulation output. These are crucial properties to understand the kinetic bottlenecks of a system, and in **Chapter IV** a Markov state model interpretation of the PPTIS path ensembles is introduced from which these properties can be extracted.

**1 d**: PPTIS has no implementation coupling with common MD engines

- While toy models are ideal to test the limitations and strengths of path sampling methodologies, their application to real-world biophysical systems ultimately decides their utility. The PPTIS and REPPTIS methodologies are implemented in the third installment of the PyRETIS software package, as detailed at the end of **Chapter III**. This research objective therefore also contains testing of the methodologies on realistic systems.

While **Chapter III** shows how REPPTIS successfully addresses the ergodic sampling issues of PPTIS for low dimensional systems, it need still be demonstrated on a challenging biophysical system. This is first done in **Chapter IV**, where REPPTIS and the new MSM analysis framework are applied to the dissociation of the drug molecule benzamidine from the trypsin protein. A protein-drug complex was chosen due to the increasing importance of virtual kinetic screening

in the early stages of the drug-design pipeline. Traditionally, equilibrium thermodynamic properties such as binding affinity (dissociation constant $K_d$) and IC$_{50}$ (concentration causing 50% target inhibition) have dominated as primary predictors of drug efficacy [158–161]. However, the importance of protein-drug kinetics, particularly the drug residence time, has been increasingly recognized for its crucial role in pharmacodynamics and better correlation with *in vivo* efficacy [162–166]. Given the costly and time-consuming nature of experimental assays, the need for virtual screening methods in the early stages of the drug-design pipeline has long been established, underscored by a growing necessity for virtual kinetic assays [162, 167–170]. This need is further emphasized by the advent of personalized medicine, where protein mutations can alter drug efficacies, requiring therapeutic strategies that are more tailored to the individual patient [171–173]. It is within this context that REPPTIS is next applied to a series of challenging protein-drug complexes in **Chapter V**. This entails a study on the dissociation of the drug molecule imatinib from the ABL kinase protein, and mutated variants of ABL.

Applications of REPPTIS on the ABL-imatinib system revealed that the methodology still faces ergodicity limitations in the presence of strongly bound metastable states. Based on the MSM interpretation of the PPTIS path ensembles, a novel path sampling methodology is proposed in **Chapter VI**, pushing towards increased ensemble communication and increased path memory.

The final part of this work delves into the study of slow oxygen kinetics within the (molecularly) large myelin sheaths. The brain covers over 20 % of the body's total oxygen metabolism, while only being 2 % of the body's weight [174]. In the brain, neurons use up to 75 % to 80 % of the brain's energy, mainly through mitochondrial oxidative phosphorylation of adenosine triphosphate [175]. Despite the critical role of oxygen, the complex pathways of oxygen from capillary to mitochondria remain largely uncharted due to experimental limitations in probing oxygen at the subcellular scale. The second main research question aims to contribute to closing this subcellar knowledge gap, by shedding light on oxygen transport through myelin sheaths.

**Research Question 2**: *How does myelination affect the storage of oxygen and its transport to axons?*

Myelin sheaths can consist of up to 100 phospholipid bilayers [176]. Simulating the transport through the entire myelin sheath is compu-

tationally infeasible due to both the large system size and the long timescales involved. Path sampling was, however, not required for this study. Instead, the periodicity of myelin sheaths was exploited to construct a diffusive model based on the simulation data of a single phospholipid bilayer, as detailed in **Chapter VII**. This model is then modified to study oxygen transport at various operating conditions and various levels of myelination.

Finally, **Chapter VIII** concludes this work and provides an outlook on future research directions.

## 1.5 LIST OF PUBLICATIONS

The research presented in this thesis has resulted in three publications (**papers I, II, and IV**), one paper that has been submitted (**paper V**), and two manuscripts that are currently in preparation (**papers III and VI**),

I   W. Vervust, D. T. Zhang, T. S. van Erp and A. Ghysels, *Path sampling with memory reduction and replica exchange to reach long permeation timescales*, Biophys. J., vol. 122, no. 14, pp. 2960-2972, 2023.

II  W. Vervust, D. T. Zhang, A. Ghysels, S. Roet, T. S. van Erp and E. Riccardi, *Pyretis 3: Conquering rare and slow events without boundaries*, J. Comput. Chem., 2024.

III W. Vervust, E. Wils and A. Ghysels, *Estimating full path lengths and kinetics from partial path transition interface sampling simulations*, **Manuscript in preparation**.

IV  W. Vervust and A. Ghysels, *Oxygen storage in stacked phospholipid membranes under an oxygen gradient as a model for myelin sheaths*, Adv. Exp. Med. Biol. 2022, 1395, 301–307

V   W. Vervust, K. Witschas, L. Leybaert and A. Ghysels, *Myelin sheaths can act as compact temporary oxygen storage units as modeled by an electrical RC circuit model*, **Submitted to PRX**.

VI  W. Vervust, E. Riccardi, D. T. Zhang, T. S. van Erp and A. Ghysels, *path sampling simulations of ABL-imatinib complexes*, **Manuscript in preparation**.

# 2

# PATH SAMPLING

This chapter provides the theoretical background of the TIS and PPTIS methodologies. First, the concepts of path space and path ensembles are introduced, also providing an overview of a particularly useful path-generating methodology for biological systems: MD simulations. Afterwards, a detailed description of the path ensembles encountered in TIS and PPTIS is given. It is then shown how rate constants can be derived from these path ensembles. Next, the core algorithm for sampling these ensembles is presented. The chapter continues by showing how incorporation of a replica exchange move has greatly improved the TIS methodology, and concludes with a 'how-to' section on performing a TIS-based simulation. Most of the content related to path sampling in this chapter (and more) can be found in the excellent book chapters of Refs. [76, 113, 177].

## 2.1 PATH SPACE AND PATH ENSEMBLES

The phase space $\mathcal{X}$ of an $M$-particle system is the $6M$-dimensional space $\mathbb{R}^{6M}$ of $x,y,z$-positions and $x,y,z$-momenta of all particles. A path $x(T)$ of duration $T$ is then defined as a sequence of $N+1$ phase points separated by a timestep $\Delta t$

$$x(T) \equiv \{x_0, x_{\Delta t}, x_{2\Delta t}, \ldots, x_{N\Delta t} = x_T\}. \tag{2.1}$$

Such a path can, for instance, be generated by an MD simulation with timestep $\Delta t$.

The path space $\mathcal{X}(T)$ is now defined as the set of all possible paths $x(T)$. The probability $\mathcal{P}[x(T)]$ of observing a particular path $x(T)$ depends on (1) the dynamics of the system (stochastic or deterministic) and (2) the initial condition of the system (phase point densities). For a Markovian process, the path probability contains a product of memoryless step transitions

$$\mathcal{P}[x(T)] = \rho(x_0) \prod_{i=0}^{N-1} p(x_{i\Delta t} \to x_{(i+1)\Delta t}), \qquad (2.2)$$

where $\rho(x_0)$ is the initial phase point density and $p(x_{i\Delta t} \to x_{(i+1)\Delta t})$ is the transition probability from phase point $x_{i\Delta t}$ to $x_{(i+1)\Delta t}$.

Path sampling methods focus their sampling on the rare regions of phase space, and so-called path ensembles are defined which concern specific subsets of path space. The simplest case is for transition path sampling, which focuses on reactive paths that start in the reactant state $A$ and end in the product state $B$. The accompanying path probability $\mathcal{P}_{AB}$ is given by

$$\mathcal{P}_{AB}[x(T)] = \frac{\theta_A(x_0)\theta_B(x_T)\mathcal{P}[x(T)]}{Z_{AB}}, \qquad (2.3)$$

where $\theta_A$ ($\theta_B$) is the indicator function that is 1 when the phase point is in state $A$ ($B$) and 0 otherwise, and where

$$Z_{AB} = \int \mathcal{D}x(T)\, \theta_A(x_0)\mathcal{P}[x(T)]\theta_B(x_T), \qquad (2.4)$$

$$= \int \cdots \int dx_0 \ldots dx_T \, \theta_A(x_0)\theta_B(x_T)\mathcal{P}[x(T)], \qquad (2.5)$$

is a normalization constant (as the total path space has been restricted). This looks very similar to configuration based ensembles, where $Z_{AB}$ is analogous to a partition function, but in path space. The above reactive path probability defines the Transition Path Ensemble (TPE), offering a complete statistical description of all possible reactive pathways.

### 2.1.1   phase point density and dynamics

In most biological simulations, the system will be put in contact with a heat bath to maintain a constant temperature $T$, as this is often used in *in vitro* experiments. The initial conditions of the system are then given by the canonical distribution

$$\rho(x_0) = \frac{\exp(-\beta H(x_0))}{Z}, \qquad (2.6)$$

where

$$Z = \int \mathrm{d}x \exp(-\beta H(x)) \tag{2.7}$$

is the partition function of the system, $\beta = 1/(k_B T)$ is the inverse temperature with $k_B$ the Boltzmann constant, and $H(x)$ is the Hamiltonian of the system.

The dynamics of MD simulations are governed by Newton's equations, and are therefore deterministic. As such, the probability of a phase point $x_t$ later on in the trajectory can be described by a propagator $\phi$ applied to the initial phase point $x_0$, where $x_t = \phi_t(x_0)$. The dynamics remain deterministic in the presence of a heat bath when, for example, the Nosé-Hoover thermostat is used. In this case, the phasepoint includes auxiliary variables $x^{\mathrm{NH}}$ that, apart from the $x$, $y$, $z$-position and velocity variables, enter the state definition $x$. Low dimensional toy systems are used to test path sampling methods, where Langevin dynamics are often used to introduce stochasticity.

The sampling strategies employed in path sampling methods often require propagation backwards in time, which for the time-reversible dynamics is equivalent to reversing the momenta and propagation forward in time. Another important property of both Langevin and Newtonian dynamics is that they obey the microscopic reversibility condition [104]

$$\rho(x)p(x \to y) = \rho(y)\bar{p}(y \to x), \tag{2.8}$$

where $\bar{p}(y \to x)$ is the time-reversed transition probability, i.e. the probability that backwards in time propagation from $y$ brings the system to $x$.

### 2.1.2 Molecular dynamics and force fields

MD simulations can provide paths by letting an initial configuration of the system evolve in time according to Newton's equations of motion. Defining $r = \{r_i\}_{i=1}^{N}$ and $p = \{p_i\}_{i=1}^{N}$ as the position- and momentum vectors of the $N$ particles/atoms contained in the system, respectively, the Hamiltonian $H$ is given by

$$H(r,p) = \sum_{i=1}^{N} \frac{p_i^2}{2m_i} + U(r), \tag{2.9}$$

where $U(r)$ is the potential energy of the system at configuration $r$. The equations of motion are given by

$$\frac{\mathrm{d}r_i}{\mathrm{d}t} = \nabla_{p_i} H = \frac{p_i}{m_i}, \tag{2.10}$$

$$\frac{\mathrm{d}p_i}{\mathrm{d}t} = -\nabla_{r_i} H = -\nabla_{r_i} U = F_i, \tag{2.11}$$

where $F_i$ is the force acting on atom $i$ with mass $m_i$. Numerical integration is then used to propagate the system in time, where at each timestep the forces acting on the atoms are calculated to update the position- and momentum vectors.

Calculation of the forces $F_i$ requires the spatial derivative of the potential energy surface $U$, which is calculated using a *force field* rather than being derived from first principles. A force field $V(r)$ approximates the potential energy surface $U(r)$ using a set of simple $r$-dependent functions. The parameters of these functions are chosen such that they optimally reproduce experimental or *ab initio* data. A force field is built up by contributions of the bonded and non-bonded interactions in the system

$$U(r) = U_{\text{bonded}}(r) + U_{\text{non-bonded}}(r). \tag{2.12}$$

The bonded interactions are designed to represent the structural characteristics of covalent bonds, including terms that approximate bond lengths, bond angles, and dihedral angles:

$$
\begin{aligned}
V_{bonded} = & \sum_{\text{bonds}} k_b (r - r_0)^2 \\
& + \sum_{\text{angles}} k_\theta (\theta - \theta_0)^2 \\
& + \sum_{\text{dihedrals}} V_n \left( 1 + cos(n\phi - \gamma) \right),
\end{aligned}
\tag{2.13}
$$

where $(k_b, r_0)$ and $(k_\theta, \theta_0)$ represent the (force constant, equilibrium position) for bond stretching between two atoms, and bond angle bending between three atoms, respectively. The third term models the torsional rotation around bonds connecting four atoms, where $(V_n, n, \gamma)$ represent the (force constant, periodicity, phase) of the torsional energy.

The non-bonded interactions describe the forces between atoms that are not directly connected by covalent bonds. These interactions include van der Waals forces and electrostatic interactions which are

crucial for capturing the behavior of condensed phases typically encountered in biological systems. The potential energy for non-bonded interactions is represented as

$$V_{\text{non-bonded}} = \sum_{i<j} \left[ \frac{A_{ij}}{r_{ij}^{12}} - \frac{B_{ij}}{r_{ij}^6} + \frac{q_i q_j}{4\pi\epsilon_0 r_{ij}} \right], \qquad (2.14)$$

where $r_{ij}$ is the distance between atoms $i$ and $j$. The first two terms capture van der Waals interactions using the Lennard-Jones potential, containing the short-range repulsive force ($\propto r^{-12}$, Pauli repulsion) and attractive forces due to induced dipoles ($\propto r^{-6}$, London dispersion). The third term describes electrostatic interactions, where $q_i$ and $q_j$ represent the partial charges of atoms $i$ and $j$, and $\epsilon_0$ is the permittivity of free space. This term captures the long-range interaction between charged atoms (Coulombic interaction).

A force field thus reduces a highly complex potential energy surface to a functional form. The art of force field development consists of the correct parameterization of all contributing interaction potentials such that MD simulations reproduce experimental and *ab initio* data. Force field validation comprises structural information (e.g. protein structures from X-ray crystallography or NMR spectroscopy data), thermodynamic properties (e.g. binding free energies, enthalpies of solvation, phase transition properties), dynamic properties (e.g. diffusion coefficients, viscosity, from NMR relaxation data or fluorescence spectroscopy), or quantum mechanical benchmarks (e.g. for small molecules). Validation is mainly based on ensemble-averaged data, which may comprise the accuracy of individual molecular events such as binding kinetics and transient conformations [178, 179]. TPS-based methodologies generate transient trajectories as modeled by MD simulations, where subsequent calculation of rate constants can be compared to experimental data. A large effort in the optimization of these methods could, in the future, result in a large enough 'transient event' dataset for force field optimization.

### 2.1.3 Transition interface sampling

Central to the TIS method is the definition of overall states $\mathcal{A}$ and $\mathcal{B}$ that partition phase space in two regions: phase points that have most recently visited $A$ or $B$, respectively. Thus, the overall state $\mathcal{A}$ contains the reactant basin $A$ and all phase points that, if you track the dynamics backwards in time, result in paths that started from basin $A$ rather than basin $B$. Similar reasoning applies to $\mathcal{B}$, and an illustration is provided in Fig. 2.1.

**Figure 2.1:** A schematic representation of the interfaces $\lambda_i$ used in TIS. The $\lambda_A$ defines the reactant basin $A$ (green), and $\lambda_B$ the product basin $B$ (blue). The interfaces $\lambda_A < \lambda_1 < \lambda_2 < \lambda_3 < \lambda_4 < \lambda_5 < \lambda_6 < \lambda_B$ track the progression of the reaction from $A$ to $B$. Phase points along the green paths belong to $\mathcal{A}$, as they track back to $A$, while phase points along the blue paths belong to $\mathcal{B}$, as they track back to $B$. The ensembles in which the paths are contained are annotated.

TIS introduces a set of $N + 1$ non-intersecting hypersurfaces (interfaces) $\lambda(x) = \lambda_i, \forall i \in \{0, \ldots, N\}$, in the region separating the reactant basin $A$ and the product basin $B$ (Fig. 2.1). The scalar function $\lambda(x)$ is the order parameter (OP) that tracks the progression of the reaction, where $\lambda_0 = \lambda_A$ and $\lambda_N = \lambda_B$ define the boundaries of the reactant basin $A$ ($[x \,|\, \lambda(x) < \lambda_A]$) and product basin $B$ ($[x \,|\, \lambda(x) > \lambda_B]$), respectively. To each of these interfaces $\lambda_i$ (except $\lambda_N$), a path ensemble $\left[i^+\right]$ is associated

$$\lambda_i \leftrightarrow \left[i^+\right], \quad \forall i \in \{0, \ldots, N-1\}\,.$$

Path ensemble $\left[i^+\right]$ contains all the paths that start from $\lambda_A$, cross interface $\lambda_i$ before recrossing $\lambda_A$, and end at either $\lambda_A$ or $\lambda_B$, as illustrated in Fig. 2.2A. For $\left[0^+\right]$, the crossing condition is automatically satisfied, as leaving state $A$ implies crossing $\lambda_0 = \lambda_A$. To implement these path ensembles computationally, a choice must be made on how to define 'starting', 'crossing', and 'ending' at an interface, due to the discrete nature of MD trajectories. Starting at $\lambda_A$ is equivalent to having the first and second phase point left and right of $\lambda_A$, respectively. Ending at $\lambda_A$ is the reverse of this, where the last and second-to-last phase point are left and right of $\lambda_A$, respectively. For a path starting and ending at $\lambda_A$ (no recrossings), the crossing condition for an ensemble $\left[i^+\right]$ equates to the path having a phase point for which $\lambda > \lambda_i$. Since the first and last phase points of a path are

**Figure 2.2:** A schematic representation of the RETIS (**A**) path ensemble $[i^+]$ and the PPTIS (**B**) path ensemble $[i^\pm]$ associated with $\lambda_i$. RETIS $[i^+]$ paths can be seen as LML or LMR paths, where the paths must start from the L(eft) interface $\lambda_A$, cross the (M)iddle interface $\lambda_i$, and end at either the L(eft) or R(ight) interfaces $\lambda_A$ or $\lambda_B$, respectively. PPTIS $[i^\pm]$ paths have their L and R boundaries at the neighboring interfaces $\lambda_{i-1}$ and $\lambda_{i+1}$, respectively. These interfaces form the start and stopping conditions for the paths, where paths must also conform to the crossing condition of $\lambda_i$. PPTIS paths can thus be any of four types: LML, LMR, RML, or RMR.

in $A$ or $B$, they are not a part of the path ensemble and are discarded in the computational analysis after a TIS simulation.

The forward rate constant $k_{AB}$ can be calculated as

$$k_{AB} = f_A P_A \left( \lambda_B | \lambda_A \right), \tag{2.15}$$

with $f_A$ the conditional flux (the flux of paths *leaving* state $A$), and $P_A \left( \lambda_B | \lambda_A \right)$ the probability that a path that has just positively crossed $\lambda_A$ will reach $\lambda_B$ before recrossing $\lambda_A$. The rate is given as the product of a flux (how many times is the process attempted?) and a crossing probability (is an attempt successful?). This global or overall crossing probability is usually extremely small, and the TIS path ensemble definitions allow this probability to be broken down as a product of history dependent interface crossing probabilities

$$k_{AB} = f_A P_A \left( \lambda_B | \lambda_A \right) = f_A \prod_{i=0}^{N-1} P_A \left( \lambda_{i+1} | \lambda_i \right), \tag{2.16}$$

where $P_A \left( \lambda_{i+1} | \lambda_i \right)$ denotes the probability that a trajectory starting from $A$ crosses the interface $\lambda_{i+1}$ after having crossed the interface $\lambda_i$ without returning first to $A$. These conditional crossing probabilities

are much larger than the overall $P_A\left(\lambda_B|\lambda_A\right)$, greatly reducing the computational cost. The conditional probabilities $P_A\left(\lambda_{i+1}|\lambda_i\right)$ are readily available from the $\left[i^+\right]$ paths, as it is simply the fraction of $\left[i^+\right]$ paths that cross $\lambda_{i+1}$. As such, each of the ensembles focuses on the calculation of one of these conditional crossing probabilities. To obtain the flux $f_A$, a conventional MD simulation is used to count the number of positive crossings of $\lambda_A$ per unit time.

In TIS, the $\lambda$ parameter is called an order parameter (OP) rather than a collective variable (CV) or reaction coordinate (RC), which calls for a short discussion (a longer discussion can be found in Ref. [117]). The combination of history dependent crossing probabilities and the importance sampling in path space (Sec. 2.2) result in a rate estimate that is exact and independent of the choice of $\lambda$, as long as it can differentiate between the reactant and product states. While the choice of $\lambda$ is, however, crucial for the convergence rate of the method, it does not affect the accuracy of the rate estimate once converged. This is in contrast with methods like TST and RF, where the rate estimate is highly dependent on the choice of $\lambda$ parameter. For these methods, deviations of the $\lambda$ parameter from the committor function (which maps configurations to their probability of hitting the product state rather than the reactant state) can result in a poor rate estimate. As such, the use of OP rather than CV or RC is used to emphasize this difference.

Sometimes, no OP exists for which TIS converges efficiently, especially for diffusive processes with many metastable states. As paths get stuck in metastable states, the paths rarely commit to states $A$ or $B$, resulting in infeasibly long paths and poor acceptance rates. Assuming that paths relax within these metastable states, it can be assumed that long time correlations become negligible and that path memory can be truncated. This is the idea behind PPTIS, where the path ensembles $\left[i^\pm\right]$ are restricted in phase space close to their $\lambda_i$ interface, resulting in interface crossing probabilities that carry less path memory.

### 2.1.4   Partial path TIS

A PPTIS path ensemble $\left[i^\pm\right]$ contains all the paths crossing $\lambda_i$ that started from $\lambda_{i-1}$ or $\lambda_{i+1}$, and ended at $\lambda_{i-1}$ or $\lambda_{i+1}$, as shown in Fig. 2.2B. As such, all the phase points of a PPTIS path ensemble are restricted to the region $[\lambda_{i-1}, \lambda_{i+1}]$. The interface crossing probabilities will be denoted by use of the generic notation $P\left(\begin{smallmatrix} C \\ D \end{smallmatrix}\middle|\begin{smallmatrix} B \\ A \end{smallmatrix}\right)$, denoting the probability that a path that had crossed $\lambda_A$ and sub-

sequently crossed $\lambda_B$ without recrossing $\lambda_A$ will cross $\lambda_C$ before crossing $\lambda_D$. The progression of such a path should thus be read $A \to B \to C$ before $D$. For each path ensemble $\left[i^{\pm}\right]$ ($i \in [1, N-1]$), local interface crossing probabilities can be calculated

$$p_i^= \equiv P\left(\begin{smallmatrix} i-1 \\ i+1 \end{smallmatrix} \middle| \begin{smallmatrix} i \\ i-1 \end{smallmatrix}\right), \qquad p_i^{\pm} \equiv P\left(\begin{smallmatrix} i+1 \\ i-1 \end{smallmatrix} \middle| \begin{smallmatrix} i \\ i-1 \end{smallmatrix}\right), \qquad p_i^= + p_i^{\pm} \equiv 1,$$

$$p_i^{\mp} \equiv P\left(\begin{smallmatrix} i-1 \\ i+1 \end{smallmatrix} \middle| \begin{smallmatrix} i \\ i+1 \end{smallmatrix}\right), \qquad p_i^{\ddagger} \equiv P\left(\begin{smallmatrix} i+1 \\ i-1 \end{smallmatrix} \middle| \begin{smallmatrix} i \\ i+1 \end{smallmatrix}\right), \qquad p_i^{\ddagger} + p_i^{\mp} \equiv 1.$$

For example, $p_i^{\pm}$ denotes the probability that an $\left[i^{\pm}\right]$ path that started from the left (from $\lambda_{i-1}$) will end on the right (to $\lambda_{i+1}$). Paths that started from the left will end either on the left or right, resulting in $p_i^{\pm}$ and $p_i^=$ summing to 1. These local crossing probabilities are accompanied by long-distance crossing probabilities

$$P_i^+ \equiv P\left(\begin{smallmatrix} i \\ 0 \end{smallmatrix} \middle| \begin{smallmatrix} 1 \\ 0 \end{smallmatrix}\right), \quad P_i^- \equiv P\left(\begin{smallmatrix} i \\ 0 \end{smallmatrix} \middle| \begin{smallmatrix} i-1 \\ i \end{smallmatrix}\right). \tag{2.17}$$

These can be calculated recursively, as $P_1^+ \equiv 1 \equiv P_1^-$, and

$$P_j^+ \approx \frac{p_{j-1}^{\pm} P_{j-1}^+}{p_{j-1}^{\pm} + p_{j-1}^= P_{j-1}^-}, \quad P_j^- \approx \frac{p_{j-1}^{\mp} P_{j-1}^-}{p_{j-1}^{\pm} + p_{j-1}^= P_{j-1}^-}. \tag{2.18}$$

The forward rate $k_{AB}$ is approximated as

$$k_{AB} \approx f_{A1} P_N^+, \tag{2.19}$$

where the effective positive flux $f_{A1}$ now denotes the positive flux through $\lambda_1$ rather than $\lambda_A$. Letting $\lambda_1$ approach $\lambda_A$ ($\lambda_1 = \lambda_A + \delta$), the flux $f_A$ becomes equivalent to the TIS flux. In this case $P_N^+ = P\left(\begin{smallmatrix} N \\ 0 \end{smallmatrix} \middle| \begin{smallmatrix} 1 \\ 0 \end{smallmatrix}\right) \xrightarrow{\delta \to 0} P\left(\begin{smallmatrix} B \\ A \end{smallmatrix} \middle| \begin{smallmatrix} 0+ \\ 0- \end{smallmatrix}\right)$ approximates the TIS global crossing probability.

## 2.2 PATH SAMPLING

The rate expressions for both TIS and PPTIS were seen to be the product of a flux term and a global crossing probability. The flux term was calculated by a separate MD simulation, while the global crossing probability was either a product (TIS) or function (PPTIS) of path ensemble conditional probabilities. These ensemble conditional probabilities are calculated by sampling the underlying path ensemble, which is performed by a MCMC simulation in path space, for each ensemble.

A first requirement is that of an initial path to seed the MCMC simulation. An initial path is often obtained by creation of a reactive trajectory by means of one of the enhanced sampling techniques mentioned in Sec. 1.2. For TIS, this reactive path is acceptable for all the path ensembles $[i^+]$, while for PPTIS this reactive path is cut to conform to the $[i^\pm]$ ensemble conditions (resulting in $N-1$ path segments for the $N-1$ ensembles).

### 2.2.1 Shooting move

Once the initial path is obtained, the MC loop is started to generate new paths. The core algorithm for path sampling has long been the *shooting* move, which is used in both TPS, TIS and PPTIS. Here, the algorithm for shooting [103] and its variant, aimless shooting [180], are considered. The algorithms are presented for variable path lengths, and with (uniformly) random selection of the shooting point [115].

Denote $x^{(o)}$ as the current path with length $L^{(o)}$, and denote $x^{(o)}(i)$ as the phase point $x^{(o)}_{i\Delta T}$. The (aimless) shooting move, for constant temperature simulations, works as follows

1. Choose a random phase point $x^{(o)}(i)$ of the old path ($i \in [1, L_o]$).

2. Modify the momenta of all the particles by drawing from the Boltzmann distribution and subsequent normalization to conserve kinetic energy (aimless shooting). Or, modify the momenta with small perturbations. Both schemes produce a modified phase point: $x^{(o)}(i) \to x^{(n)}(i)$.

3. Accept the new momenta with probability

$$\min\left[1, \exp\left(\beta\left(E^{(o)}(i) - E^{(n)}(i)\right)\right)\right],$$

   where $E(x)$ is the total system energy at $x$. If accepted, go to step 4, otherwise reject the move and return to step 1. For aimless shooting, this step results in automatic acceptance.

4. Draw a (uniform) number $\alpha$ from $[0, 1]$. Set the maximum path length of the new path to $L^{(n)}_{\max} = \text{int}(N^{(o)}/\alpha)$.

5. Integrate the equations of motion (EoMs) backward in time from $x^{(n)}(i)$ until either the max path length is reached (rejection, return to step 1) or until TIS: $\lambda_A$ or $\lambda_B$ is hit, or until PPTIS: $\lambda_{i-1}$ or $\lambda_{i+1}$ is hit.

6. Propagate the EoMs forward in time from $x^{(n)}(i)$ until either the max path length is reached (rejection, return to step 1) or until TIS: $\lambda_A$ or $\lambda_B$ is hit, or until PPTIS: $\lambda_{i-1}$ or $\lambda_{i+1}$ is hit. Concatenate the backward and forward segments to create the new path $x^{(n)}$. If this path is not included in the path ensemble definition, reject and return to step 1.

7. The new path is accepted, and $x^{(n)}$ becomes the new $x^{(o)}$.

While it is very intuitive to reject paths that do not conform to the path ensemble definition, it may at first glance seem like wasted compute time to reject based on path lengths and kinetic energy. However, this is required to obey detailed balance, which is of crucial importance in sampling paths according to their underlying path ensemble probabilities (i.e. their weights).

### 2.2.2   Detailed balance

To obey detailed balance, the probability of moving from one path to another must equal the probability of moving in the opposite way

$$\mathcal{P}(x^{(o)})\pi[x^{(o)} \rightarrow x^{(n)}] = \mathcal{P}(x^{(n)})\pi[x^{(n)} \rightarrow x^{(o)}], \qquad (2.20)$$

where $\mathcal{P}(x)$ denotes the path ensemble probability of path $x$ and $\pi[x \rightarrow y]$ denotes the probability of moving from path $x$ to path $y$. Obeying detailed balance is a sufficient condition to sample paths according to their respective weights (given ergodic sampling). In the Metropolis-Hastings importance sampling scheme [15], the transition probability is broken down into a generation probability $P_{\text{gen}}$ and an acceptance probability $P_{\text{acc}}$. Defining the acceptance probabilities as follows then satisfies detailed balance

$$P_{\text{acc}}(x^{(o)} \rightarrow x^{(n)}) = \theta_{[i]}(x^{(n)}) \min \left[ 1, \frac{\mathcal{P}(x^{(n)})P_{\text{gen}}\left(x^{(n)} \rightarrow x^{(o)}\right)}{\mathcal{P}(x^{(o)})P_{\text{gen}}\left(x^{(o)} \rightarrow x^{(n)}\right)}, \right]$$

$$(2.21)$$

where $\theta_{[i]}(x^{(n)})$ is 1 when $x^{(n)}$ is in the path ensemble $[i]$ ($\left[i^+\right]$ for TIS or $\left[i^{\pm}\right]$ for PPTIS), and 0 otherwise. The acceptance/rejection checks in the shooting move are thus a direct consequence of this sampling scheme.

### 2.3   REPLICA EXCHANGE TIS

Separate path sampling simulations are required in all the TIS ensembles. It is likely that the MCMC random walks of the ensembles

will explore different regions of phase space, and it is quite unfortunate that this information is not shared over all of the ensembles. Especially on the rough free energy landscape of biological systems, multiple reaction pathways will need to be explored by all the ensembles. If the exploration is hindered by slow dynamics or rather large energy barriers orthogonal to the OP, it is possible that the path ensembles sample only a few of the existing pathways, resulting in ergodic sampling issues.

The replica exchange move aims to overcome this inefficiency, by allowing ensembles to swap their paths with each other. If a swap move between ensemble $[i^+]$ and $[j^+]$ is attempted, it is simply checked whether the paths conform to the criteria of the other ensemble. If not, the swap move is rejected, and if so, the ensembles exchange paths. As was previously noted, the $[(i+1)^+]$ ensemble is a subset of the $[i^+]$ ensemble, and one only needs to check one path for acceptance.



**Figure 2.3:** A schematic representation of the replica exchange moves in RETIS. **A**: the $[i^+]$ and $[(i+1)^+]$ paths can be swapped, as both paths are contained in one another's path ensemble. **B**: the $[0^-]$ and $[0^+]$ paths can always be 'swapped', where backwards extension of the $[0^+]$ path results in a $[0^-]$ path, and forwards extension of the $[0^-]$ path results in a $[0^+]$ path.

RETIS also introduced a new ensemble $[0^-]$ near the reactant state $A$. The $[0^-]$ ensemble is defined by paths that start at $\lambda_A$, sample the reactant state $A$ $(x|\lambda(x) < \lambda_A)$, and end at $\lambda_A$. The $[0^-]$ paths exclusively sample the left side of $\lambda_A$, and can therefore never be swapped with the $[i^+]$ $(i \in [0, N-1])$ paths that sample the right side of $\lambda_A$. However, a special $[0^-] \leftrightarrow [0^+]$ swap move was

designed to allow communication between the reactant state and the transition region, as illustrated in Fig. 2.3B. This move, which is always acceptable, consists of propagating the $[0^-]$ path forward in time until it crosses either $\lambda_A$ or $\lambda_B$, resulting in a new $[0^+]$ path (after cutting away the $\lambda < \lambda_A$ part). Conversely, the $[0^+]$ path is propagated backwards in time until it hits $\lambda_A$, resulting in a $[0^-]$ path (after cutting away the $\lambda > \lambda_A$ part).

The need for a separate MD simulation to calculate $f_A$ has been alleviated, as it can be derived directly from the average path lengths of the $[0^+]$ and $[0^-]$ paths

$$ f_A = \frac{1}{\langle\tau\rangle_{[0^-]} + \langle\tau\rangle_{[0^+]}}, \tag{2.22} $$

where the brackets $\langle\tau\rangle_{[\cdot]}$ denotes a path ensemble average of the path duration $\tau$ (i.e. the product of path length and time step $\Delta t$).

The replica exchange move greatly increases the convergence rate of sampling the TIS path ensembles, where path space exploration of individual paths is now shared between the ensembles. While initially suggested for PPTIS in Ref. [121] in 2004, the replica exchange move was introduced only to TIS (Fig. 2.3), in 2007 [115]

## 2.4 Performing a RE(PP)TIS simulation

Prior to running a TIS-based simulation, some initialization is required:

1. creating initial paths for all ensembles,

2. defining an order parameter $\lambda$,

3. defining (and updating) interface positions $\{\lambda_i\}_{i=0}^N$.

Each of the (PP)TIS ensembles require an initial path to which the MC moves can be applied. Typically, a path sampling simulation is preceded by a a long MD simulation to ensure that the system is at (or near) equilibrium. A segment of this trajectory can then be cut for the $[0^-]$ ensemble. A steered MD simulation is (typically) performed to generate a reactive trajectory that conforms to all positive TIS ensembles $[i^+]$. For a PPTIS simulation, the reactive trajectory is to be cut into path segments that each conform to their respective

local ensemble definitions $[i^{\pm}]$. The initial trajectories are not required to be representative of a spontaneous reaction at equilibrium, as the Metropolis-Hastings machinery ensures that the trajectories will (eventually) be sampled according to the equilibrium path ensemble distribution. As it is expected that the first $N_{\text{init}}$ paths carry some bias due to the initialization, these are to be discarded in post-analysis (typically $N_{\text{init}} \approx 1000$ to $10\,000$). It is therefore best practice to keep the bias force in the steered MD simulation small, such that equilibration to representative path sampling happens quickly.

While the choice of order parameter $\lambda$ does not affect the accuracy of TIS simulations, it impacts the efficiency of TIS simulations (i.e. the speed of convergence). Therefore, $\lambda$ should be chosen such that it adequately captures the important barriers encountered during a reactive event. In other words, the presence of barriers orthogonal to $\lambda$ that are important to the reaction mechanism should be avoided, as the importance sampling scheme provides efficient sampling only *along* $\lambda$. This also applies to PPTIS simulations, where additionally the choice of $\lambda$ will impact the accuracy. Due to the loss of memory in PPTIS paths, so-called *tunneling effects* can occur where barriers along or orthogonal to $\lambda$ are missed. A more detailed explanation of tunneling is given in Chapter 6 and Ref. [113].

As a rule of thumb, optimal TIS efficiency requires the ensemble local crossing probabilities to be $\sim 20\,\%$ [133]. An educated guess for interface positioning can be made an estimation of the free energy profile along $\lambda$ is available, e.g. by running an umbrella simulation prior to path sampling. After a few 1000 cycles of MC path sampling moves, the initial estimates of the local crossing probabilities will indicate where along $\lambda$ the interface distribution should be more dense (if the crossing probabilities are too low) or more sparse (if the crossing probabilities are too large). The last paths of the ensembles using the initial interface placement can be re-used (sometimes requiring manipulation to obey the new ensemble crossing conditions), such that little simulation time is lost. Once interface positioning is satisfactory, the (RE)(PP)TIS simulation can run until the block averages for (1) all the local crossing probabilities and (2) the rate estimate has converged. More information on the interpretation of (RE)(PP)TIS simulation output is available at `https://pyretis.org`.

# 3

# REPPTIS: PUSHING PATH SAMPLING TO LONGER TIMESCALES

This chapter introduces the REPPTIS methodology, where inclusion of a replica exchange move (or 'swapping move') to the PPTIS framework aims to improve its accuracy. This is mainly achieved by an increase in ergodic sampling, where the swapping move allows exploration of phase space regions that are otherwise hardly accessible. The performance of the REPPTIS methodology and its implementation in PyRETIS 3 are found in **paper I** and **paper II**, respectively. This chapter aims to provide guiding and additional information to these papers, primarily through technical details and extended discussions.

This chapter therefore contributes to **Research Objectives 1a** (ergodic sampling) and **1d** (coupling with MD engines) of this work. The chapter begins by introducing the swapping move and detailing its algorithmic implementation. This is followed by a closer examination of the special case ensembles near the reactant state. The chapter concludes with a discussion on the performance of the REPPTIS method.

## 3.1 REPLICA EXCHANGE PPTIS

As discussed in chapter 2, RETIS path ensembles are subspaces of one another ($\left[(i+1)^+\right] \subset \left[i^+\right]$), where the swapping of paths required no MD integration. An exception to this was the swap between $\left[0^-\right]$ and

$[0^+]$ paths, where extremal phase points had to be extended into each other's path ensembles. The swap move proposed for PPTIS happens in a similar fashion, and is illustrated in Fig. 3.1.
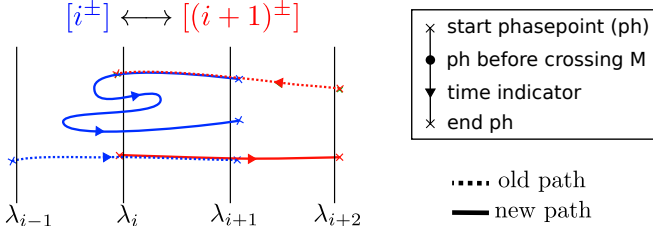


**Figure 3.1:** A schematic representation of the replica exchange move in PPTIS. The $[i^\pm]$ and $[(i+1)^\pm]$ paths can be swapped, as both paths can be extended into each other's path ensemble.

A swap move can only be performed between adjacent path ensembles $[i^\pm]$ and $[(i+1)^\pm]$. If paths of adjacent ensembles are both located in their overlapping region of phase space, then extensions of these paths (in the correct time direction) are guaranteed to produce viable paths for the other ensemble. In the example of Fig. 3.1, the last part of the RML path of $[(i+1)^\pm]$ (i.e. the *ML part in the overlapping region $[\lambda_i, \lambda_{i+1}]$) can be seen as an RM* part in $[(i+1)^\pm]$, and extension of this path is guaranteed to result in either an RML or RMR path in $[i^\pm]$ (in the example it resulted in an RMR path).

The swap move facilitates exploration of phase space in regions that are normally hardly accessible. The improved sampling strategy results in an improved accuracy in the presence of orthogonal barriers, as shown in **paper I** for a maze model system for membrane permeation. It is also shown how REPPTIS can efficiently extract the permeability for a model of ibuprofen permeation through a lipid bilayer, where the presence of multiple metastable states hinder efficient application of RETIS.

The swap move, as implemented in PyRETIS 3, is detailed in Algorithm 1. It also applies to the $[0^-]$ and $[0^{\pm\prime}]$ ensembles at the reactant state, which are discussed below. The choices of propagation directions in steps 2 and 3 are made randomly and with equal probability, which is key to maintain detailed balance. A previous implementation, used in **paper I**, only proposed forward propagation for the $[i^\pm]$ ensemble, and only backward propagation for the $[(i+1)^\pm]$ ensemble. Such implementation only probes one possible swappable configuration, and performs it whenever possible. The new implementation allows for all swappable configurations to be probed,

---

**Algorithm 1** Replica exchange move in PPTIS (swapping move)

---

1: Choose a (uniform) random number $r \in [1, \ldots, N]$. Denoting $\{E_i\}_{i=1}^{N} = \left[0^-\right], [0^{\pm\prime}], \left[1^{\pm}\right], \ldots, \left[(N-2)^{\pm}\right]$, then ensemble $E_r$ will attempt to swap paths with ensemble $E_{r+1}$.

2: Choose a random propagation direction for the old $E_r$ path $x_r^{(o)}$ and check whether this extends the path into $E_{r+1}$. If not, reject the move.

3: Choose a random propagation direction for the old $E_{r+1}$ path $x_{r+1}^{(o)}$ and check whether this extends the path into $E_r$. If not, reject the move.

4: Extend both paths in the chosen propagation directions until a stopping condition is met, resulting in new (extended) paths $x_{r,\text{ext}}^{(o)}$ and $x_{r+1,\text{ext}}^{(o)}$.

5: Remove the phase points of the extended paths that are not included in the path ensemble to which they will be exchanged, resulting in the new paths $x_{r+1}^{(n)}$ and $x_r^{(n)}$.

6: Swap these paths between the ensembles and accept the move.

---

and can be interpreted as the possibility to do a time reversal move on the paths before *and* after the swap move is performed.

## 3.2 ENSEMBLES NEAR THE REACTANT STATE

The ensembles $\left[0^-\right]$ and $[0^{\pm\prime}]$ denote the ensembles left and right of $\lambda_A$, respectively, where $\left[0^-\right]$ is identical to its definition in RETIS. The ensemble $[0^{\pm\prime}]$ is special, in the sense that it only has a left boundary at $\lambda_A$ and a right boundary at $\lambda_1$. Similarly to TIS, one could define a crossing interface $\lambda_0 + \delta$ for the $[0^{\pm\prime}]$ ensemble, where $\delta$ is an infinitesimal positive number. For paths starting from the left, the crossing condition is then automatically satisfied. These paths are therefore denoted as LL and LR paths, rather than the usual LML and LMR notation, to emphasize the lack of a true (M)iddle interface. For paths starting from the right, crossing the interface $\lambda_0 + \delta$ is practically equivalent to crossing $\lambda_0$. This means that there are only RL paths, and no RR paths, where again the M is omitted.

The swap move between $\left[0^-\right]$ and $[0^{\pm\prime}]$ is very similar to the RETIS $\left[0^-\right] \leftrightarrow \left[0^+\right]$ swap move, where now the paths can be extended in both directions. The swap move between $[0^{\pm\prime}]$ and $\left[1^{\pm}\right]$ is also possible. Paths of $\left[1^+\right]$ starting or ending at $\lambda_0$ (LMR, LML, or RML paths) have a part that can be cut to form a path in $[0^{\pm\prime}]$. On the

other hand, LR and RL paths of $[0^{\pm\prime}]$ can be extended to form LMR and RML paths in $\left[i^{\pm}\right]$.

## 3.3   CONCLUSION

The REPPTIS methodology successfully increased ergodic sampling for the model systems of **paper I**, addressing **Research Objective 1a**. The methodology was also successfully implemented in PyR-ETIS 3 in **paper II**, allowing the application of REPPTIS to complex biological systems in the following two chapters, and addressing **Research Objective 1d**. For now, extraction of the rate still requires a separate MD simulation to calculate the flux term. An expression for the flux can, however, be derived directly from the simulation output, as well as other important time-dependent quantities, which is the topic of the next chapter.

# REPPTIS AS A MARKOV STATE MODEL

The PPTIS and REPPTIS methodologies produce short paths, from which long distance crossing probabilities are calculated. A formalism to extract the temporal aspect of these longer paths is, however, lacking. In this chapter, a Markov state model formalism is introduced to tackle exactly this issue. The contents of this chapter are based on **paper III** (**manuscript in preparation**), where this chapter aims to mainly introduce the formalism through some examples and discussions that are not included in the paper.

As such, this chapter contributes to the third **Research Objective 1c** of Research Question 1, providing the framework from which time-dependent properties of long continuous paths can be extracted. The REPPTIS methodology is also applied for the first time to a biological system using MD, where the dissociation kinetics of the benzamidine molecule from the trypsin protein is investigated. As such, this chapter also contributes to the fourth **Research Objective 1d** of Research Question 1, by shedding light on the performance of REPPTIS for biological systems.

## 4.1 TRAJECTORIES AS SEQUENCES OF PPTIS PATH SEGMENTS

Consider the long MD trajectory $U$ of Fig. 4.1, shown in black. This trajectory is seen as a sequence of overlapping path segments ($U_0$,

$U_1$, $U_2$, $U_3$, ...), where each path segment $U_n$ belongs to a specific path type of a PPTIS path ensemble. These are the colored path segments of Fig. 4.1, where $U_0 \in \mathrm{LR}_{[0\pm']}$, $U_1 \in \mathrm{LMR}_{[1\pm]}$, $U_2 \in \mathrm{LML}_{[2\pm]}$, $U_3 \in \mathrm{RML}_{[1\pm]}$, etc. The four path types of $[i^\pm]$ ensembles are now
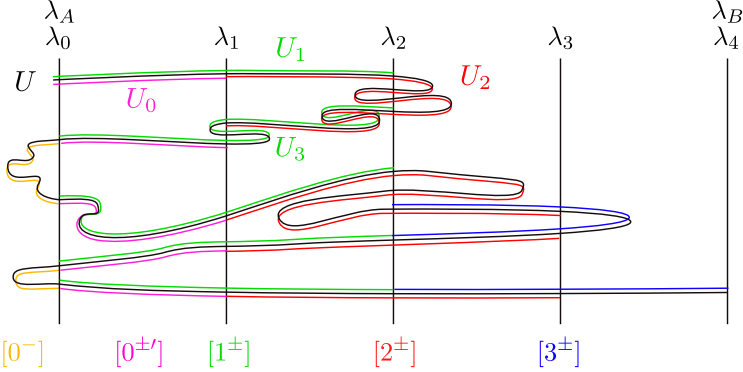


**Figure 4.1:** A long MD trajectory $U$ (black) decomposed into a sequence of path segments $(U_0, U_1, U_2, U_3, ...)$ belonging to PPTIS path ensembles. The path segments are colored according to the ensemble they belong to.

viewed as states $S_i^{k,l}$, where $k$ and $l$ denote the starting and ending interfaces of the path type, respectively. The values of $k$ and $l$ lie in $\{-1, +1\}$, where $-1$ refers to L and $+1$ refers to R. As such, $S_i^{-1,+1}$ refers to $\mathrm{LMR}_{[i\pm]}$, the LMR path types in ensemble $[i^\pm]$. Similarly, $S_i^{-1,-1}$ refers to $\mathrm{LML}_{[i\pm]}$, $S_i^{+1,-1}$ to $\mathrm{RML}_{[i\pm]}$, and $S_i^{+1,+1}$ to $\mathrm{RMR}_{[i\pm]}$. The $[0^{\pm'}]$ and $[0^-]$ ensembles are special cases, and do not contain all four path types. The $[0^{\pm'}]$ ensemble does not have a middle interface, and contains the three states $S_0^{-1,+1}$ (LR paths), $S_0^{+1,-1}$ (RL paths), and $S_0^{-1,-1}$ (LL paths), as RR paths are not sampled. The $[0^-]$ ensemble has neither left nor middle interface, and only contains the state $S_{0-}^{+1,+1}$ (RR paths). As discussed in Chapter 2, middle interfaces for the $[0^{\pm'}]$ and $[0^-]$ ensembles can be defined to lie infinitesimally close to $\lambda_A$, but the 'M' is dropped in the path type notation for simplicity. The new notation uses the numerical indices $\{k, l\}$ rather than the textual indices $\{L,R\}$ to simplify the MSM equations that follow shortly. The long MD trajectory $U = (U_0, U_1, U_2, U_3, ...)$ of Fig. 4.1 can now be mapped to the state chain $(S_0^{-1,+1}, S_1^{-1,+1}, S_2^{-1,-1}, S_1^{+1,-1}, ...)$.

## 4.2 MARKOV STATE MODEL

The state chain of the previous section can be viewed as a Markov chain through state space, where the transition probabilities between

states are defined by the local crossing probabilities of the PPTIS path ensembles. Consider, for example, an LMR path in ensemble $[2^{\pm}]$, corresponding to the red path of state $S_2^{-1,+1}$ in Fig. 4.2. If
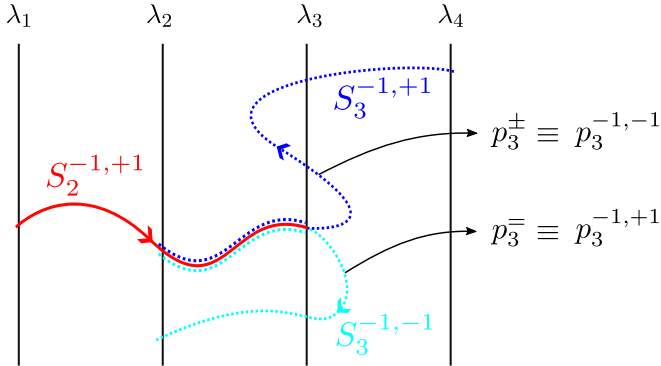


**Figure 4.2:** A path of $S_2^{-1,+1}$ (red) can either transition to state $S_3^{-1,-1,}$ (cyan) or $S_3^{-1,+1,}$ (blue). The probabilities of these transitions are given by the local crossing probabilities $p_3^{-1,-1}$ and $p_3^{-1,+1}$, respectively.

this path is propagated forwards in time until it crosses an interface *other* than $\lambda_3$, it must either end at $\lambda_2$ or $\lambda_4$. This extended path is therefore seen as one of two possible chains, $(S_2^{-1,+1}, S_3^{-1,-1})$ (cyan extension) or $(S_2^{-1,+1}, S_3^{-1,+1})$ (blue extension). The probabilities of these extensions can be estimated from the PPTIS simulation output as follows. Looking at the 'second part' of the red path, it is known that this path went from $\lambda_2$ to $\lambda_3$, corresponding to the L and M interfaces of the $[3^+]$ ensemble, respectively. The probability that this path will continue to cross the R or L interface of the $[3^{\pm}]$ ensemble is exactly the definition of the local crossing probabilities $p_3^{\pm}$ and $p_3^{=}$, respectively. For the MSM formalism, these probabilities are written with the $\{k, l\}$ notation, where now

$$
\begin{aligned}
p_i^{\pm} &= p_i^{-1,+1}, & p_i^{=} &= p_i^{-1,-1}, & p_i^{-1,-1} + p_i^{-1,+1} &= 1 \\
p_i^{\mp} &= p_i^{+1,-1}, & p_i^{\ddagger} &= p_i^{+1,+1}, & p_i^{+1,-1} + p_i^{+1,+1} &= 1.
\end{aligned}
\tag{4.1}
$$

As such, it is seen that a state can only transition to two other states, the probability of which is given by the relevant local crossing probabilities. This is concisely captured in a transition matrix $M$

$$
M_{ikl,i'k'l'} = P\left(S_i^{k,l} \to S_{i'}^{k',l'}\right),
\tag{4.2}
$$

with elements

$$
M_{ikl,i'k'l'} = \begin{cases} p_{i+l}^{-l,+1}, & i' = i + l, k' = -l, l' = +1 \\ p_{i+l}^{-l,-1}, & i' = i + l, k' = -l, l' = -1 \\ 0, & \text{elsewhere} \end{cases}
\tag{4.3}
$$

While the use of indices is quite heavy, they indicate that the end point $l$ of the segment $S_i^{k,l}$ determines whether the next state is in a higher ensemble $i' = i + 1$ ($l = +1$, indicating the end point R) or in a lower ensemble $i' = i - 1$ ($l = -1$, indicating the end point L). This clarifies that the next visited path ensemble will be $\left[i'^{\pm}\right]$ with $i' = i + l$. Furthermore, the end point $l$ determines the starting point of the next segment, making the next starting point $k' = -l$. The new segment can then have an arbitrary end point, so $l'$ can be either $+1$ or -1. In **paper III**, the three-label indices $ikl$ are flattened to a single Greek index $\alpha$, to alleviate the notational load of the equations that follow.

## 4.3   MSM analysis

The statistics of long paths can now be analyzed using MSM theory, where the transition matrix $M$ can be used to calculate *hitting* and *return* probabilities that allow for estimation of long distance crossing probabilities. Each state, considering it represents a specific path type, has an associated average path length $\tau$. To calculate the length of larger paths, the *accumulated* time of walks through the MSM is considered, which is the sum of the path lengths of the states visited. Special care is however required to correctly account for the overlap of PPTIS path segments, such that the accumulated time is not overestimated. To this end, path lengths $\tau$ associated to $S_i^{k,l}$ states are decomposed into three parts (Fig. 4.3): the part $\tau_{(1)}$ before the first crossing of $\lambda_i$, the part $\tau_{(2)}$ after the last crossing of $\lambda_i$, and the part $\tau_{(m)}$ in between. Visiting a state then contributes $\tau_{(m2)} = \tau_{(m)} + \tau_{(2)}$ to the accumulated time of the walk.



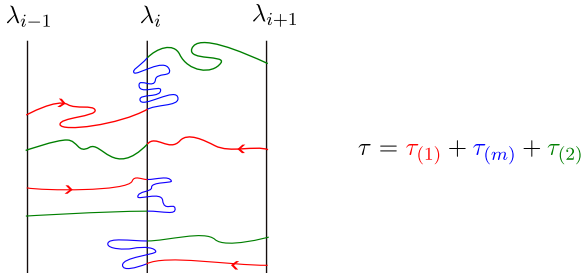$$\tau = \tau_{(1)} + \tau_{(m)} + \tau_{(2)}$$

**Figure 4.3:** The path lengths $\tau$ of $\left[i^{\pm}\right]$ states are decomposed into three parts: $\tau_{(1)}$ before the first crossing of $\lambda_i$, $\tau_{(2)}$ after the last crossing of $\lambda_i$, and $\tau_{(m)}$ in between. These parts can be zero, as seen for the second path without a (m)iddle part.

Of particular interest are the average path lengths $\tau_{[0^-]}$ and $\tau_{[0^+]}$. Their sum determines the conditional flux $f_A = (\tau_{[0^-]} + \tau_{[0^+]})^{-1}$ (Eq. 2.22) and subsequently the reaction rate $k_{AB} = f_A P_A(A \to B)$ (Eq. **??**). In **paper III**, it is shown how these quantities are calculated with the MSM formalism, after which the methodology is verified with low-dimensional model systems by direct comparison with RETIS results. It is currently still being explored how an $A \to B$ MFPT distribution can be extracted from the MSM framework. As all possible $A \to B$ path probabilities and their respective lengths can be calculated, the MSM formalism could allow for kinetic determination of reactions for which the rate is ill-defined (e.g. reactions that are not well-described by a single rate constant).

## 4.4 Trypsin-benzamidine dissociation

**Paper III** also holds the application of REPPTIS to study the dissociation kinetics of the benzamidine molecule from the trypsin protein, a biological system that is often used to benchmark enhanced sampling methodologies. REPPTIS had difficulty in sampling the $\left[1^{\pm}\right], \left[2^{\pm}\right], \ldots, \left[6^{\pm}\right]$ ensembles close to the reactant state, where the free energy profile is steepest. It was assumed that, due to a poor initialization of the paths in these ensembles, the MC moves were unable to relax towards more favorable regions of path space.

A long MD simulation was used to estimate crossing probabilities associated to these ensembles, whose incorporation then resulted in a rate constant that was in good agreement with experimental data and other computational work. As unbinding of benzamidine is characterized by an initial steep energy profile, the contribution of $\tau_{[0^+]}$ to the flux is minimal. Therefore, the flux estimate using the MSM formalism was in good agreement with the MD simulation, despite poor sampling of the first few ensembles.

The replica exchange moves were infrequently performed in ensembles close to the bound state, where it was hypothesized that a combination of local orthogonal barriers and local steep free energy profiles caused ergodic sampling issues. This is because the steep energy profile results in dominant sampling of LML paths (rolling down the barrier). If this behavior spans over multiple neighboring ensembles, paths are rarely exchanged, as an $*MR\leftrightarrow LM*$ connection is required (where $*$ can be either L or R). The shooting move, on the other hand, enables efficient sampling mainly *along* the $\lambda$ parameter, and not orthogonal to it. As the initial path was created with a

steered MD simulation containing a directional bias, it is possible that benzamidine did not escape along a dominant reactive pathway, and instead sampled a metastable region whose separation from the dominant pathway was not well captured by the $\lambda$ parameter. While it was expected that the replica exchange move would have resolved such initialization issues for the benzamidine-trypsin system, it showed that correct path ensemble initialization is crucial for REPPTIS performance.

To conclude this chapter, the MSM formalism was introduced as a tool to analyze the statistics of long paths generated by the short (RE)PPTIS paths, allowing for the extraction of the flux and rate constant from the simulation output. As such, **Research Objective 1c** was addressed. The application to the benzamidine-trypsin system was in line with **Research Objective 1d**, where the performance of REPPTIS to larger biological systems was explored.

# ABL IMATINIB DISSOCIATION

This chapter contributes to **Research Objective 1d** of Research Question 1, by taking an in-depth look at the performance and limitations of the REPPTIS methodology to complex biological systems.

Traditionally, equilibrium thermodynamic properties such as binding affinity (dissociation constant $K_d$) and IC$_{50}$ (concentration causing 50% target inhibition) have dominated as primary predictors of drug efficacy [158–161]. However, the importance of protein-drug kinetics, particularly the drug residence time, has been increasingly recognized for its crucial role in pharmacodynamics and better correlation with *in vivo* efficacy [162–166]. Given the costly and time-consuming nature of experimental assays, the need for virtual screening methods in the early stages of the drug-design pipeline has long been established, underscored by a growing necessity for virtual kinetic assays [162, 167–170]. This need is further emphasized by the advent of personalized medicine, where protein mutations can alter drug efficacies, requiring therapeutic strategies that are more tailored to the individual patient [171–173].

## 5.1 INTRODUCTION

The Abelson nonreceptor tyrosine kinase (ABL) plays a significant role in signal transduction, cell differentiation, cell division, cell adhesion and stress response [181, 182]. The ABL kinase domain catalyzes the transfer of the $\gamma$-phosphate from adenosine triphosphate

(ATP) to tyrosine residues in substrate proteins [183]. Patients with the fusion oncogene *BCR-ABL* suffer from chronic myeloid leukemia (CML), where the normally autoinhibited kinase domain becomes constitutively active [182, 184]. The tyrosine kinase inhibitor (TKI) imatinib (Gleevec®, also known as STI571) was developed to competitively bind to the ATP-binding site of ABL, which was a hallmark for molecular targeted therapies in cancer [182, 184, 185]. Mutations in the ABL kinase domain resulted in imatinib resistance, to which second and third generation TKIs were developed [184, 186].



**Figure 5.1:** Inactive (**A**, PDB 6NPV [187]) and active (**B**, PDB 6XR6 [188]) conformations of the ABL kinase domain. **C**: Close up of the binding pocket in the inactive state. Imatinib binds to the inactive state, where the aspartate residue points outwards of the binding pocket. **D**: Close up of the binding pocket in the active state. In the active state, the aspartate residue points inwards of the binding pocket. In **A** and **B**, the dashed cyan lines represent missing residues of the 6NPV structure. In **D**, the dashed cyan lines represent removed residues for better visualization of the binding pocket. Figure created with PyMOL [189].

The ABL kinase domain is shown in Fig. 5.1, consisting of the N-terminal and C-terminal lobes typical for protein kinases. Three important regions are highlighted: the activation loop (A-loop, red),

the ATP-binding loop (P-loop, green), and the $\alpha$C-helix (magenta). The N-lobe consists of a five stranded $\beta$-sheet and the $\alpha$C-helix, and also contains the P-loop. The C-lobe mainly consists of helical structures and contains the substrate binding site [183]. The active site (ATP-binding site) is located in the cleft between the N- and C-lobes.

Located at the N-terminal side of the A-loop is the aspartate-phenylalanine-glycine (DFG) motif. The active state is characterized by the aspartate pointing inwards towards the ATP-binding site (DFG-in, Fig. 5.1D). Transitioning to the inactive state is accompanied by the DFG motif flipping $180°$, where the aspartate residue moves outwards and the phenylalanine moves inwards (DFG-out, Fig. 5.1C) [183]. In the active state, the $\alpha$C-helix is turned away from the hinge region, and the A-loop is in an open (extended) conformation. In the inactive state, the $\alpha$C-helix is turned towards the hinge region, and the A-loop can block substrate binding to the C-lobe.

In this study, the dissociation of imatinib from the wild-type (WT) ABL kinase domain is investigated using the REPPTIS method. Furthermore, 7 mutated variants of ABL are studied (Table 5.1 and Fig. 5.2A). To run REPPTIS simulations, initial paths for the ensembles are required, which are obtained by first running equilibrium simulations followed by steered MD simulations.

## 5.2 PREPARATION FOR REPPTIS SIMULATIONS

### 5.2.1 Equilibrium simulations

The crystal structure of the ABL kinase domain in complex with imatinib was obtained from the 6NPV PDB entry [187]. As there were

| Variant | $H_2O$ | $Na^+$ | $Cl^-$ |
|---------|--------|--------|--------|
| WT | 17355 | 61 | 53 |
| Q252H | 17359 | 61 | 53 |
| Y253F | 17360 | 61 | 53 |
| E255V | 17362 | 60 | 53 |
| T315I | 18858 | 62 | 56 |
| M351T | 17358 | 61 | 53 |
| F359V | 17364 | 61 | 53 |
| H396P | 17358 | 61 | 53 |

**Table 5.1:** Number of $H_2O$ molecules, $Na^+$ ions, and $Cl^-$ ions used in the (mutated) ABL systems. WT denotes the wild type ABL kinase

no structures available for 6 of the 7 mutated ABL variants in complex with imatinib or similar TKIs, CHARMM-GUI [190] was used to mutate the residues based on the 6NPV structure. For the gatekeeper mutation T315I, a crystal structure was available in complex with the inhibitor DCC-2036 (PDB entry 3QRJ [191]). Both the 6NPV (blue) and 3QRJ (green) ABL structures are shown in Fig. 5.2B. The 3QRJ structure was then superimposed on the 6NPV structure, with the DCC-2036 ligand removed and the imatinib ligand added.



**Figure 5.2:** **A**: Wild type ABL (gray structure) with the 7 mutations considered. **B**: The ABL kinase domain in complex with imatinib (PDB 6NPV, blue), and in complex with DCC-2036 (PDB 3QRJ, red). Figure created with VMD [192].

These structures were then prepared for MD simulations using the Gromacs [37] software package (version 2021.3). The CHARMM36m force field [193] was used for the protein and solvent, where the TIP3P water model was used [194]. The force field parameters for imatinib were taken from Ref. [195], which were obtained using QM calculations using the GAAMP software [196, 197] and which were recently used in another study on the WT ABL-imatinib complex [198]. The ligand force field was provided in CHARMM format, and it was converted to the GROMACS format using CHARMM GUI [190]. The imatinib parameters define a protonated state, which is desirable as at physiological pH the complex with neutral imatinib represents less than 0.1% of the overall population [199].

The structures were solvated in a rhombic dodecahedron box, ensuring a minimum distance of 1.5 nm between periodic images. The dodecahedral box is only 71% of the cubic box volume with the same

minimum distance criterion, resulting in significant computational savings. The minimal image distance is chosen to be larger than the typical 1 nm value as dissociation is to be studied. $Na^+$ and $CL^-$ ions were added to neutralize the system and to reach a physiological salt concentration of 0.15 M. The water and ion content of the simulation boxes is shown in Table 5.1. On average, the 8 systems contained 17547 waters, 61 $Na^+$ ions, and 53 $Cl^-$ ions. The systems were then energy minimized using steepest descent until the maximum force was below $1000 \, \mathrm{kJ \, mol^{-1} \, nm^{-1}}$. The systems were then equilibrated in the NVT and NPT ensemble for 100 ps, each, where both the ligand and protein heavy atoms were restrained. The production simulations (unrestrained) were run for 50 ns each, in the NPT ensemble at 300 K and at 1 bar using the Bussi–Donadio–Parrinello thermostat [200] (coupling constant of 0.1 ps) and the Parinello-Rahman barostat [201] (coupling constant of 2 ps and compressibility of $4.5 \times 10^{-5} \, \mathrm{bar^{-1}}$). Intramolecular hydrogen bonds with heavy atoms were constrained using LINCS [52], which allowed for a time step of 2 fs.

The production simulations were checked for validity (temperature, energy, pressure) and by visualizing the trajectories with the VMD [192] software. The root-mean-square deviation (RMSD) and root mean square fluctuation (RMSF) of the protein $C_\alpha$'s are shown in Fig. 5.3, where protein alignment was performed using the *rigid* $C_\alpha$'s only (defined in Sec. 5.2.2.1), where the first frame served as reference.

### 5.2.2 Generating an initial path

#### 5.2.2.1 The order parameter

The order parameter $\lambda(x)$ for a configuration (time slice, frame) $x$ is defined as the distance between the current center of mass (COM) of imatinib $\mathbf{r}_{\mathrm{COM}}(x)$ and the average bound COM position of imatinib $\mathbf{r}_{\mathrm{COM}}^{\mathrm{bound}}$.

$$\lambda(t) = \left\| \mathbf{r}_{\mathrm{COM}}(t) - \mathbf{r}_{\mathrm{COM}}^{\mathrm{bound}} \right\| \tag{5.1}$$

where $\mathbf{r}_{\mathrm{COM}}^{\mathrm{bound}}$ is defined as follows.

The $C_\alpha$ atoms of the residues with an RMSF value lower than a threshold value $\mathrm{RMSF}_{max}$ are defined as rigid residues (see Table 5.2 for the values used per variant). Residues belonging to the A-loop, the P-loop, or the $\alpha$C-helix were excluded from being defined as rigid residues. For each variant production trajectory, the frames were superposed on the *final* frame (i.e. the reference frame $x_{\mathrm{ref}}$) using the $C_\alpha$'s of the rigid residues, after which the average COM position of
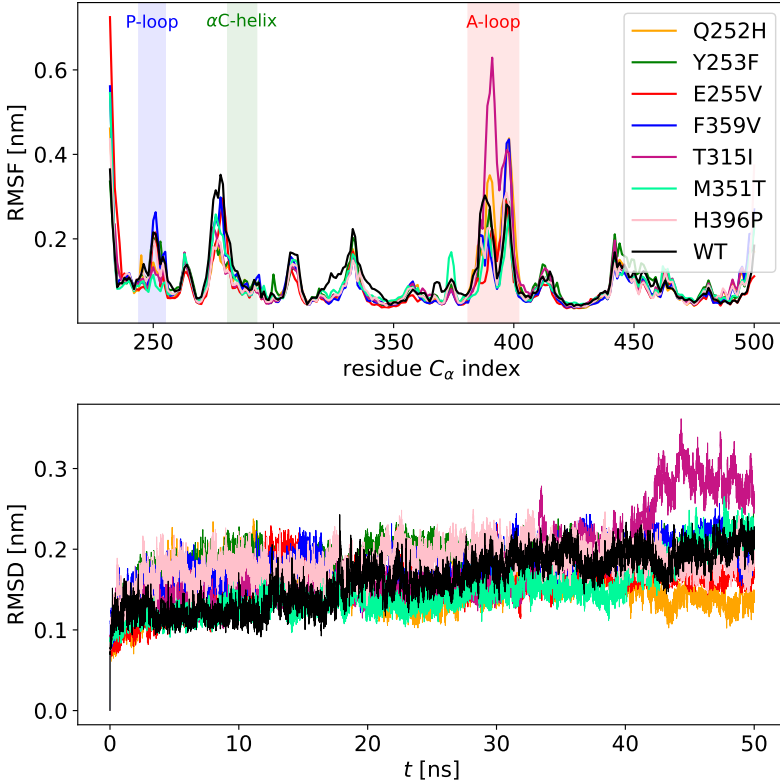
**Figure 5.3:** Production run RMSF (**A**) and RMSD (**B**) of the $C_\alpha$ atoms, for native ABL (WT, black line) and the 7 mutated variants (colored lines). The first frame of each production run served as reference, and the alignment was performed using the *rigid* $C_\alpha$'s (as defined in Sec. 5.2.2.1). Both figures share the same legend. Residues are numbered according to the UniprotKB accession number P00519 (tyrosine-protein kinase ABL1).

imatinib was calculated and defined as $\mathbf{r}_{\mathrm{COM}}^{\mathrm{bound}}$. The order parameter $\lambda(x)$ of a new frame $x$ can then be calculated by first superposing this frame to the reference frame (using the rigid $C_\alpha$'s), after which the current COM position of imatinib is calculated, and its distance to $\mathbf{r}_{\mathrm{COM}}^{\mathrm{bound}}$ is calculated. The order parameters calculated for the production trajectories are shown in Fig. 5.4, where metastable states are seen to be present even in the deep binding pocket of ABL. This is particularly visible for the Q252H variant (Fig. 5.4A), where the first 12 ns of the trajectory are spent in a metastable state that is clearly separated from the remaining 38 ns. It was chosen not to incorporate the initial 12 ns of the Q252H production run for the definition of its average bound COM position $\mathbf{r}_{\mathrm{COM}}^{\mathrm{bound}}$.
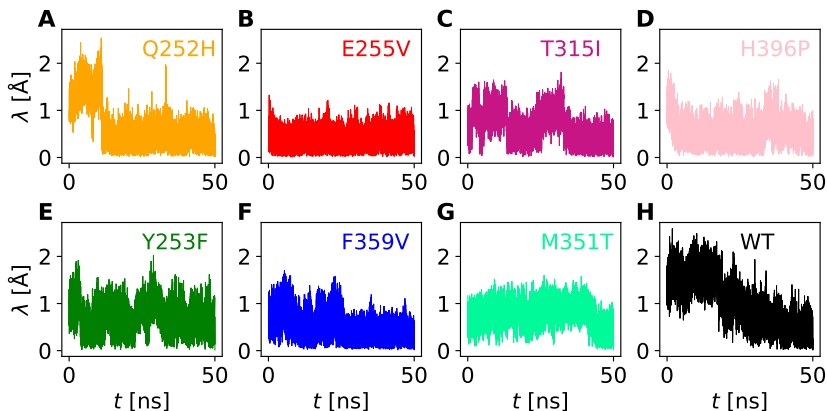
**Figure 5.4:** The order parameter $\lambda$ trajectories for the production runs of the ABL variants.

| Variant | Rigid residues |
|---------|----------------|
| WT | 276, 316-321, 333-337, 339-343, 345, 359-372, 424-426, 438-454, 488-496, 502-504, 508 |
| Q252H | 316-318, 332-337, 341-343, 358-373, 424, 426, 439-455, 458, 488, 490-495, 504-512 |
| Y253F | 316-318, 320-321, 332-334, 336-337, 341-342, 358-374, 379, 426, 437-450, 453, 490-496, 502-506, 508-512 |
| E255V | 316-318, 320-322, 333-337, 341-346, 358-372, 380, 422-426, 437-456, 458, 488-496, 502-508 |
| T315I | 317-318, 320-323, 332-337, 342, 361-372, 423-426, 438-454, 458, 476-477, 490-496, 502, 504-510, 512 |
| M351T | 316-318, 320-321, 332-337, 341-343, 360, 362-373, 423-424, 438-453, 458, 477, 488-496, 502-503, 508 |
| F359V | 268-270, 295-302, 313-318, 320-324, 339-354, 364, 368-372, 376-380, 418-435, 457-458, 471-478, 483-486, 489 |
| H396P | 314, 316-317, 320-323, 332-337, 365-374, 424, 426, 437-454, 487-496, 502-513 |

**Table 5.2:** The residues defined as rigid residues for the calculation of the order parameter $\lambda$. Numbering according to the UniprotKB accession number P00519 (tyrosine-protein kinase ABL1).

#### 5.2.2.2 Pulling simulations

The pulling simulations (steered MD) were performed with Gromacs using the PLUMED [202] extension (version 2.7.2). Care was taken to ensure that the pulling simulations did not introduce large artifacts in the dissociation pathways, which involved manually checking

the trajectories and several reruns with optimized parameters. Artifacts were encountered in two ways. First, it was noticed that pulling without restraints on the protein resulted in the C-lobe and N-lobe separating (the hinge region completely opening up). Second it was noticed that t pulling force and velocity should be chosen carefully. If these values are too high, the system does not have time to relax the bias force into the DoFs orthogonal to $\lambda$ (relax to local equilibrium), which can greatly impact the equilibration time of REPPTIS (as discussed later on).

Good pulling trajectories were obtained by restraining the $C_\alpha$'s of the rigid residues ($\kappa = 5000 \, \text{kJ} \, \text{mol}^{-1} \, \text{nm}^{-2}$) and by pulling along the order parameter $\lambda$ with a moving restraint ($\kappa = 250 \, \text{kJ} \, \text{mol}^{-1} \, \text{nm}^{-2}$) aimed to displace imatinib 38 Å from the bound position in 50 ns (i.e. a pulling velocity of $0.76 \, \text{Å} \, \text{ns}^{-1}$). As the pulling was done slowly, and $\lambda$ is defined as the distance to the average bound state, it is expected that the pulling force was non-directional. Six of the variants escaped via a pathway under the $\alpha$C-helix, while the Q252H and F359V variants had imatinib escape via a pathway under the P-loop. The potential energy of the systems during the pulling simulations are shown in Fig. 5.5A, where the perturbation introduced by the pulling force is hardly visible. The evolution of the order parameter $\lambda$ during pulling is shown in Figs. 5.5B-C, where the characteristic nature of rare event transitions, marked by sudden jumps in the $\lambda$ parameter, is seen to be preserved (Fig. 5.5C). If the pulling force is too high, this plot would contain straight lines, where dissociation is more likely to sample unlikely pathways. Imatinib is strongly bound in the binding pocket of ABL, as suggested by a sudden transition (jump in $\lambda$) around the 5 ns to 10 ns mark, depending on the variant. Most of the variants do not simply display a single transition, where for example a double transition for WT ABL is clearly visible. This suggests the presence of metastable states, some of which have imatinib still strongly bound to ABL. Once imatinib transitions out of the deep binding pocket and its metastable states, the pulling profile is rather linear. This suggests that, at least for these pathways, the ligand was no longer strongly bound to the protein.

### 5.2.2.3   Umbrella simulation

The shape of the free energy profile for WT ABL was probed with an umbrella simulation (Fig. 5.6). This was performed on an implementation of the system with a smaller simulation box (minimal image distance of 1 nm) and where ions where placed only to neutralize the system (not to reach a physiological salt concentration). The only objective of the umbrella simulation is to obtain a first guess for the interface

**Figure 5.5:** **A**: The potential energy of the systems during the pulling simulations, void of large artifacts. The purple line represents the T315I variant, which was constructed from the 3QRJ PDB entry and has a larger simulation box (and hence lower potential energy). **B**: Evolution of the order parameter $\lambda$ during the pulling simulations. **C**: Zoom-in on the first 15 ns of dissociation, showing 'jumping behavior' typical of activated processes. Native ABL (black line) shows a double transition, indicating the presence of a strongly bound metastable state.

placements, and it was therefore not rerun for the larger simulation box used in the REPPTIS simulations. The umbrella simulations were performed with Gromacs using 10 ns runs for 52 windows, totalling

520 ns. Of these 52 windows, 28 used $\kappa = 1000\,\mathrm{kJ\,mol^{-1}\,nm^{-2}}$ and the additional 24 windows were run with $\kappa = 5000\,\mathrm{kJ\,mol^{-1}\,nm^{-2}}$ at lower $\lambda$ values for better coverage (Fig. 5.6A). The free energy profile (Fig. 5.6B) was reconstructed with Gromacs' built-in WHAM approach, where standard deviations are provided by bootstrapping (200 resampled datasets with replacement) [203].

The steep increase of the free energy well is indicative for imatinib being strongly bound within the binding pocket, where ~80 % of the barrier is overcome within the first 10 Å. As such, the distance between interfaces $\lambda_i$ needs to be small for the first 10 Å, while it can be larger for larger $\lambda$ values.



**Figure 5.6:** Umbrella simulation on the (native) ABL-imatinib complex. **A**: Histograms of the 52 umbrella windows, where the 28 red windows used a force constant $k = 1000\,\mathrm{kJ\,mol^{-1}\,nm^{-1}}$, and the 24 extra windows used a force constant $k = 5000\,\mathrm{kJ\,mol^{-1}\,nm^{-1}}$ (black). **B**: The reconstructed free energy profile of the WT ABL-imatinib complex using all 52 windows. The shade represents standard deviations obtained by bootstrapping, where the data was resampled 200 times with replacement.

### 5.2.3 Trajectory-wise order parameter calculation

The $\lambda$ OP is not calculated for each MD timestep (2 fs) as this would impede MD simulation rates require infeasible storage space. Instead, $\lambda$ is calculated each $n_{\text{subcycles}}$ timesteps. In the original Gromacs 1 engine of PyRETIS, the MD integration and $\lambda$ calculation were performed iteratively, where Gromacs is launched to perform $n_{\text{subcycles}}$ MD steps, shut down, $\lambda$ calculated to check if a stopping criterion is met, and if not, this process is repeated. The Gromacs 2 engine drastically increased the efficiency by running MD integration and $\lambda$ calculation in parallel. However, the extreme efficiency of MD software introduced a new bottleneck, where the $\lambda$ calculation can severely lag behind the MD engine. While this lag is rather small for simple order parameters (e.g. atom-atom distances using the internal PyRETIS TRR reader), it becomes significantly limiting for more complex $\lambda$ definitions. Loading a frame from a trajectory using MDTraj [204], for example, can take up to 1 s, which is often a lot longer than the time to perform $n_{\text{subcycles}}$ MD steps. As such, MD integration can reach its maximal allowed length, while the $\lambda$ calculation has only reached 10% of the frames. As such, compute time is wasted on 'dead' time, where no MD steps are performed. MDAnalysis [205] is better suited for consecutively loading frames, as it allows predefining a 'Universe' object for which the topology is loaded *a priori*, after which frames can be loaded quickly ($\ll 1$ s).

Calculating $\lambda$ for the ABL-imatinib simulations requires:

1. loading frame $x(t)$ (cheap with MDAnalysis),

2. making the protein and ligand whole, i.e. no broken bonds over PBCs (expensive),

3. Aligning $x(t)$ onto $x_{\text{ref}}$ using rigid $C_\alpha$'s (expensive),

4. calculating the imatinib $\text{COM}(t)$-$\text{COM}_{\text{ref}}$ distance (cheap, but tricky).

The second and third steps are the most expensive, where a frame-by-frame calculation of $\lambda$ cannot be made competitive with the MD engine. This bottleneck can be alleviated by calculating $\lambda$ for subtrajectories rather than single frames. This is because calculation (using existing optimized algorithms for steps 2 and 3) of *the next frame* in a trajectory is a lot cheaper than calculation of *a single frame* trajectory. The $\lambda$ calculation was changed to fit this new scheme. Denote

$t_{\mathrm{MD}}$ as the current length of the MD trajectory and $t_{\mathrm{OP}}$ as the frame for which $\lambda$ is to be calculated. By construction $t_{\mathrm{OP}} \leq t_{\mathrm{MD}}$, and one of two things happens in the $\lambda$ calculation loop:

1. Look up whether $\lambda(t_{\mathrm{OP}})$ was already calculated and stored in the list $L_\lambda$ of precalculated values. If it is, return this value. If it is not, go to step 2.

2. Calculate $\lambda$ for the subtrajectory interval $[t_{\mathrm{OP}}, t_{\mathrm{MD}}]$, and save to a list $L_\lambda = [\lambda(t_{\mathrm{OP}}), \lambda(t_{\mathrm{OP}} + 1), \ldots, \lambda(t_{\mathrm{MD}})]$.

Using this scheme, the rate of OP calculation for the ABL-imatinib systems increased approximately 50-fold, practically nullifying 'dead' compute time.

Apart from the need for competitive $\lambda$ calculation speeds, care is required to avoid any PBC artifacts. As a rhombic dodecahedron simulation box in the NPT ensemble is used, the triclinic unit cell vectors are not constant. The safest way to calculate $\lambda$ is to superpose the reference protein-ligand complex onto that of the current frame, such that *(un)wrapping* behavior is based on the unit cell vectors of the current frame. Distances are then also to be calculated w.r.t. the unit cell vectors of the current frame.

The trajectory-based OP calculation was (successfully) tested for both reproducibility and comparison with the original frame-based OP calculation. The eventual algorithm used a combination of MDAnalysis functions and GROMACS *trjconv* calls.

### 5.2.4 $\infty$REPPTIS

Further large speed-up of the REPPTIS simulations was achieved by implementing REPPTIS in the $\infty$RETIS software, with the help of Daniel T. Zhang (NTNU). The $\infty$RETIS [118, 119] software allows multiple *workers* to run MD moves (or other MD intensive moves) in parallel. This implementation was intended solely for the ability to use more hardware resources, as the infinite swapping formalism is not applicable to REPPTIS ensembles. The software was (successfully) tested for usage with the internal PyRETIS engine and the external Gromacs engine. The $\lambda$ calculation for the ABL-imatinib systems, however, had to be reimplemented as the MDAnalysis functions were not compatible with the $\infty$RETIS software. The new OP implementation was based on the PLUMED plugin using the *FIT_-TO_TEMPLATE* function [202]. PLUMED couples to most common

MD engines, and is often used for enhanced sampling methods or analysis methods that require CV calculations. As PLUMED has direct access to the internal coordinates of the MD engine, its on-the-fly CV calculations are fast, making PLUMED a promising candidate for future path sampling software. For the ABL-imatinib system, the $\lambda$ calculation is performed every $n_{\text{subcycles}} = 20$ MD steps by the Gromacs-PLUMED run itself. The calculated $\lambda$ values are stored to a list that is then just read in by the PyRETIS $\lambda$ calculator.

## 5.3 REPPTIS SIMULATIONS

Initially, RETIS was used to study the ABL-imatinib systems. Nearly all shooting moves (and later wire-fencing moves) resulted in trajectories that did not commit to $\lambda_A$ or $\lambda_B$ within the maximum allowed path length, even for values as high as 50 ns. This was attributed to the presence of long-lived metastable states along the dissociation pathways, where the extremely low acceptance rate made RETIS infeasible. REPPTIS was then used to study imatinib dissociation, for which the key input parameters are shown in Fig. 5.7. Order parameters were calculated every 40 fs, and the maximum path length was set to 4 ns. Paths exceeding this length are forcefully rejected to restrict disk space usage, and their occurrence should be minimal as they break detailed balance. Occurrences were rare, as they only happened [0, 0, 1, 3, 3, 25, 0, 7] times for the [WT, Q252H, Y253F, E255V, F359V, T315I, M351T, H396P] variants. For all variants, 45 interfaces were used, where the first interfaces are closely spaced (deep binding pocket) and then gradually spaced further apart (Fig. 5.7). The steered MD simulations of Sec. 5.2.2.2 were used to initialize the REPPTIS path ensembles. Swapping moves were attempted for 25% of the MC moves, and the remaining 75% were shooting moves.

```
Simulation
----------
task = repptis
interfaces = [1, 1.25, 1.5, 1.75, 2, 2.25, 2.5, 2.75, 3,   # every 0.25 Å
              3.33, 3.66, 4, 4.33, 4.66, 5, 5.33, 5.66, 6,  # every 0.33 Å
              6.5, 7, 7.5, 8, 8.5, 9, 9.5, 10, 10.5, 11,    # every 0.5 Å
              12, 13, 14, 15, 16, 17, 18, 19, 20            # every 1 Å
              22, 24, 26, 28, 30, 32, 34, 36, 38]           # every 2 Å

Engine settings
---------------
class = gromacs3      # sub-trajectory OP calculation†
timestep = 0.002      # MD integration timestep 2 fs
subcycles = 20        # OP calculated every 40 fs

TIS settings
------------
maxlength = 100000   # maximum path length 4 ns
aimless = True       # aimless shooting algorithm

RETIS settings
--------------
swapfreq = 0.25      # 25% swap moves, 75% shooting moves
```

**Figure 5.7:** Selection of key parameters used in the REPPTIS ABL-imatinib simulations (in *.rst* format). Interfaces are placed close to each other in the deep binding pocket of ABL, after which they are gradually spaced further apart. The subcycles parameter ($n_{\mathrm{subcycles}} = 20$) denotes the number of MD steps between OP calculations (i.e. REPPTIS 'sees' every 20th frame). †: Adaptations to the gromacs2 engine were made to run the subtrajectory OP calculation.

| variant | $N_{\mathrm{MC}}$ | $N_{\mathrm{ACC}}$ | problematic ensembles | $P_A(B\|A)$ $[10^{-27}]$ | $f_A$ $[\mathrm{ps}^{-1}]$ | $k_{\mathrm{AB}}$ $[\mathrm{s}^{-1}]$ |
|---|---|---|---|---|---|---|
| WT | 53967 | 15003 | $[0-15]$, 19, 25, 29, 36, $[40-41]$ | † | $4.2 \cdot 10^{-2}$ $\pm 19\%$ | † |
| Q252H | 7156 | 2499 | $[0-45]$ | † | † | † |
| Y53F | 45642 | 19402 | $[0-16]$, 38 | † | $2.4 \cdot 10^{-1}$ $\pm 9.3\%$ | † |
| E255V | 72283 | 32597 | $[0-10]$, 30, 31 | 1.3 $\pm 70\%$ | $1.6 \cdot 10^{-2}$ $\pm 30\%$ | $2.1 \cdot 10^{-17}$ $\pm 76\%$ |
| T315I | 108733 | 46370 | $[0-15]$, $[34-36]$ | 32 $\pm 53\%$ | $9.6 \cdot 10^{-3}$ $\pm 21\%$ | $3.1 \cdot 10^{-16}$ $\pm 57\%$ |
| M351T | 44451 | 21481 | $[0-6]$, 17, 20, 21, 27 | † | $9.8 \cdot 10^{-1}$ $\pm 24\%$ | † |
| F359V | 17157 | 6006 | $[0-45]$ | † | $1.6 \cdot 10^{-3}$ $\pm 85\%$ | † |
| H396P | 62497 | 17694 | $[0-4]$, 18, 19, 32, 33, $[40-42]$ | † | $2.6 \cdot 10^{-3}$ $\pm 31\%$ | † |

| variant | method | $k_{\mathrm{off}}[\mathrm{s}^{-1}]$ |
|---|---|---|
| WT [206] | experiment | $(2.2 \pm 4.6) \cdot 10^{-3}$ |
| WT [36] | experiment | $(8.3 \pm 0.8) \cdot 10^{-4}$ |
| WT [207] | experiment | $(2.5 \pm 0.6) \cdot 10$ |
| WT [208] | milestoning | $1.8 \cdot 10$ |
| WT [209] | InMetaD‡ | $(6 \pm 3) \cdot 10^{-4}$ |

**Table 5.3:** Results of the ABL-imatinib REPPTIS simulations (top half of the table) and rate estimates obtained from literature (bottom half of the table). None of the simulations have converged to a (reliable) rate constant. All simulations used 46 ensembles. $N_{\mathrm{MC}}$ and $N_{\mathrm{ACC}}$ denote the amount of MC moves performed and the amount of MC moves accepted, respectively. Problematic ensembles are discussed in the text.
†: not enough data was gathered to calculate the global crossing probability $P_A(B|A)$, and/or the rate constant $k_{AB}$, and/or the flux $f_A$.
‡: infrequent metadynamics (InMetaD) simulation.

The results of these simulations are shown in Table 5.3, where the initial path was excluded from the analysis. None of the simulations converged to a reliable rate estimate, where most did not contain sufficient data to calculate the global crossing probability $P_A(\lambda_B|\lambda_A)$ and the rate constant $k_{AB}$. While a rate constant and error estimate was obtained for the E255V and T315I variants, their values are not reliable. due to failed sampling in a large portion of the path ensembles ('problematic ensembles' column of Table 5.3), which is discussed in more detail in the next section. These ensembles failed to sample representative paths, resulting in extremely small (or zero-valued) local

crossing probabilities. To get an idea of the long crossing probability profiles $P_i^+ = P_A(\lambda_i|\lambda_A)$, zero-valued local crossing probabilities were artificially set to 0.01, where the complementary local crossing probabilities were then set to 0.99 to keep their sum at 1 (e.g. $p_i^{\pm} + p_i^{=} \equiv 1$). The resulting profiles are shown in Fig. 5.8. While a large drop in crossing probability is expected to escape the deep binding pocket, the magnitude of the drop is largely overestimated.



**Figure 5.8:** The crossing probability profiles $P_A(\lambda_i|\lambda_A)$ for the ABL-imatinib systems over the entire $\lambda \in [1, 38]\,\text{Å}$ range (**A**) and zoomed in on the $\lambda \in [1, 6]\,\text{Å}$ range (**B**). Profiles related to the REPPTIS simulations (lines with X-markers for each interface) have artificially assigned local crossing probabilities in case of zeros (see text). The black dashed line denotes the $\infty$RETIS simulation (see Sec. 5.3.2). The black dash-dotted line (WT$^{\dagger}$) is constructed by appending the $\lambda \in [6, 38]\,\text{Å}$ profile of WT REPPTIS to the $\lambda \in [1, 6]\,\text{Å}$ profile of WT $\infty$RETIS (see Sec. 5.3.2).

### 5.3.1 REPPTIS sampling issues

First, a note on the rate error estimations is due. The error estimates for the E255V and T315I rates are not reliable due to a severe under-estimation of the global crossing probability $P_A(\lambda_B|\lambda_A)$ errors. This is due to the (recursive) block averaging performed on the running estimate of $P_A(\lambda_B|\lambda_A)$, which is only justified if the errors for the local crossing probabilities are (at least) of the same order of magnitude. This is not the case for the ABL-imatinib systems, where the problematic ensembles sampled only very few paths of certain types (e.g. LMR) in the entire simulation. As such, the error estimate for the corresponding local crossing probabilities (e.g. $p_i^{\pm}$) are large, which are lost in the block analysis on the running estimate of $P_A(\lambda_B|\lambda_A)$. This is not to be considered a deep flaw of the error estimation for $P_A(\lambda_B|\lambda_A)$ in general, as it should only be used *if* its constituent local crossing probabilities are well sampled (i.e. the sampling converges). To conclude, the *precision* of the rate estimates was overestimated.

The ABL-imatinib systems did not converge in the problematic ensembles. For all the simulations, these include the ensembles close to the deep binding pocket (small $\lambda$ values), and some ensembles further along the unbinding pathway. This was not due to a severe lack of accepted paths in these ensembles, as there is no large discrepancy between the amount of accepted paths $N_{\mathrm{ACC}}(i)$ for the different ensembles (Fig. 5.9).

Two issues are assumed to be at the root of the sampling issues. First, the one-dimensional order parameter $\lambda$ is not able to differentiate the slow dynamics well enough. That means that, for some ensembles, energy barriers along $\lambda^{\perp}$ (i.e. the degrees of freedom orthogonal $\lambda$) can be important or even dominant for the dissociation process. This then couples with the second main issue, where the initial path may not have been representative of a true unbiased pathway. The orthogonal barriers then hinder the path sampling procedure to explore the relevant regions of phase space. As an example, it may be possible that a rotation of a dihedral angle between imatinib rings is the rate limiting step close to the deep binding pocket. If this rotation is missed in the initial path, then it is unlikely to be introduced by the sampling procedure if this rotational energy barrier is large.

The average path lengths $\langle\tau\rangle_i \forall i \in [1, 46]$ of the ensembles are also shown in Fig. 5.9, for each of the variants. The path length is measured in 'number of phasepoints', where the starting and ending phase points of paths are excluded. The path lengths are much lower near the binding pocket, which is to be expected as the energy barrier

is steepest in this region. As the occurrence of LMR or RMR paths is rare in these ensembles, the REPPTIS swapping moves designed to improve phase space exploration were rarely performed.
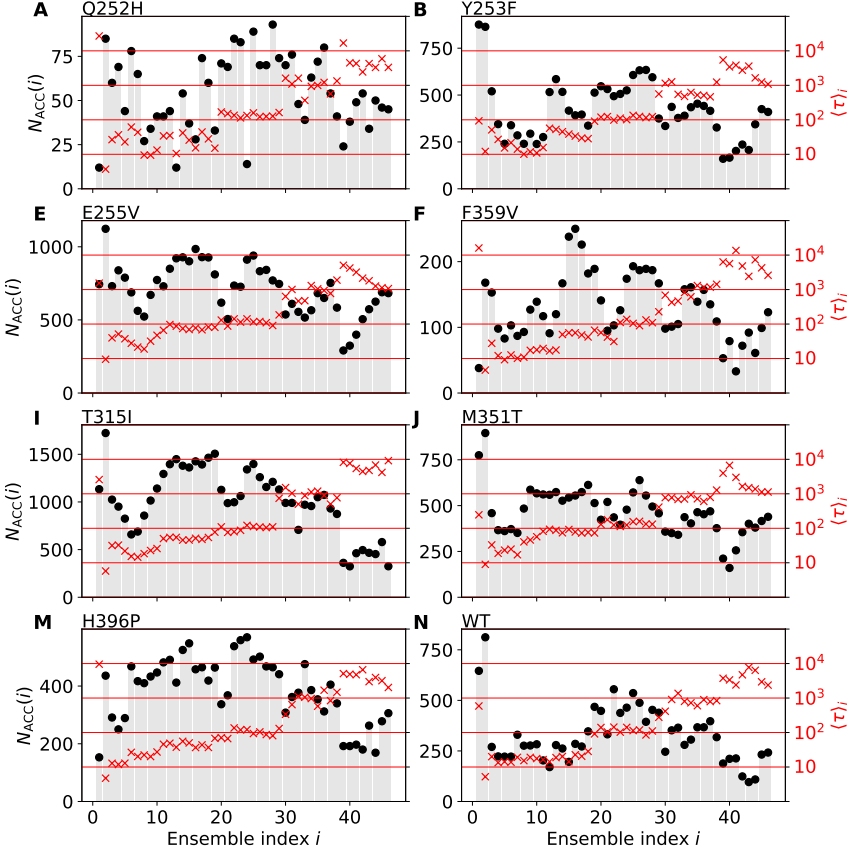


**Figure 5.9:** The amount of unique accepted paths $N_{\mathrm{ACC}}(i), \forall i \in [1,46]$ for each of the 46 ensembles in the REPPTIS simulations is shown with a gray bar plot and black dots, for each of the variants. Ensemble index 1 corresponds to $[0^-]$, index 2 to $[0^{\pm\prime}]$, index 3 to $[1^\pm]$, etc. Note that each variant uses a different scale for $N_{\mathrm{ACC}}(i)$ (left y-axis). The average path lengths $\langle\tau\rangle_i$ of the ensembles are shown with red X-marks, for each of the variants. This represents the average number of phasepoints (excluding end and start point) of $[i^\pm]$ paths. Note that each variant uses the same logarithmic scale for $\langle\tau_i\rangle$ (right y-axis).

Let us now focus on WT ABL, whose crossing probability profile $P_A(\lambda_i|\lambda_A)$ showed the fastest drop of all the variants (Fig. 5.8). Path lengths for the WT ABL system in the $\left[[1^+], \ldots, [16^+]\right]$ ensembles (ensemble indices 3 to 18 on Fig. 5.9N) vary around 20, with

approximately 250 accepted paths per ensemble. This corresponds to $20 \times 40\,\mathrm{fs} \times 250 = 200\,\mathrm{ps}$ of combined trajectory lengths, for each of these ensembles. Only a fraction of this time can be seen as 'equilibration time', as the shooting points are (on average) not located near the beginning or end of paths. It is expected that the initial path of WT ABL was not representative for the dominant reaction pathway, where the steered MD simulation pulled it along a barrier orthogonal to $\lambda$. This barrier was then not overcome by the fraction of the 200 ps of phase space exploration. Therefore, the local crossing probabilities $p_i^{\pm}$ in these ensembles remain extremely low (or zero). For problematic ensembles further along the unbinding pathway ($\left[i^{\pm}\right]$, with $i \geq 20$), the path lengths are considerably longer. However, the efficiency of phase space (and path space) exploration is not solely determined by the path lengths, but more so by the height of the orthogonal energy barriers. It is expected that, for these ensembles, orthogonal barriers are also present, hindering efficient sampling of the corresponding $\left[i^{\pm}\right]$ path spaces. The above reasoning also applies to the other variants.

To further investigate the sampling issues close to the binding pocket, an $\infty$RETIS simulation was run for the WT ABL system for the $\lambda \in [1, 6]\,\text{Å}$ range.

### 5.3.2  $\infty$**RETIS**

The input parameters used for the $\infty$RETIS simulation of WT ABL are shown in Fig. 5.10. Interfaces were placed from $\lambda_A = 1\,\text{Å}$ to $\lambda_B = 6\,\text{Å}$ with a separation of 0.1 Å, resulting in 51 interfaces and 51 ensembles $\{E_i\}_{i=0}^{50} = \left\{\left[0^-\right], \left[0^+\right], \left[1^+\right], \ldots, \left[49^+\right]\right\}$. The interfaces can be put in closer proximity to each other, as the infinite swapping formalism of $\infty$RETIS efficiently distributes (swaps) path information between the ensembles. Furthermore, the optimal amount of workers (MD intensive moves that can run in parallel) is approximately half of the number of ensembles, where it is thus advantageous to place more interfaces than a PyRETIS RETIS simulation.

The resulting crossing probability profile $P_A(\lambda|\lambda_A)$ is included in Fig. 5.8A-B (black dashed line) for comparison with the REPPTIS profiles, and also shown separately in Fig. 5.11A. The first 5000 accepted paths were discarded in the analysis to avoid initialization effects. The WHAM approach of RETIS produces a continuous profile, rather than discrete points at the interface positions for REPPTIS. The drop in crossing probability, while less steep than the WT REPPTIS profile, is still exaggerated, especially in the $\lambda \in [1, 2]\,\text{Å}$ range. Dis-

```
[dask]
workers = 20  # amount of workers executing MD intensive moves

[simulation]
interfaces = [1, 1.1, 1.2, 1.3, 1.4, 1.5, 1.6, 1.7, 1.8, 1.9, # every 0.1Å
              2, 2.1, 2.2, 2.3, 2.4, 2.5, 2.6, 2.7, 2.8, 2.9, # every 0.1Å
              3, 3.1, 3.2, 3.3, 3.4, 3.5, 3.6, 3.7, 3.8, 3.9, # every 0.1Å
              4, 4.1, 4.2, 4.3, 4.4, 4.5, 4.6, 4.7, 4.8, 4.9, # every 0.1Å
              5, 5.1, 5.2, 5.3, 5.4, 5.5, 5.6, 5.7, 5.8, 5.9] # every 0.1Å
              6]                                               # every 0.1Å
shooting_moves = ['sh']×51   # shooting moves in all ensembles

[engine]
class = gromacs      # this is equivalent to PyRETIS gromacs2†
timestep = 0.002     # MD integration timestep 2 fs
subcycles = 500      # OP calculated every 1 ps
gmx_format = 'g96'   # higher precision for dense interfaces

[simulation.tis_set]
maxlength = 100000  # maximum path length 100 ns
aimless = true      # aimless shooting algorithm
```

**Figure 5.10:** Selection of key parameters used in WT $\infty$RETIS simulation (in *.toml* format). The subcycles parameter ($n_{\text{subcycles}} = 500$) denotes the number of MD steps between OP calculations (i.e. $\infty$RETIS 'sees' every 500th frame). $^{\dagger}$: Adaptations to the gromacs engine were made to run the subtrajectory OP calculation.

continuities in the $P_A(\lambda|\lambda_A)$ profile are visible (arrow indications on Fig. 5.11A), and are discussed next.

The discontinuity around $\lambda \approx 1.9\,\text{Å}$ is due to failed sampling in the first 9 positive ensembles $[0^+], \ldots, [8^+]$. The average paths lengths of these ensembles (now measured in ps, as $n_{\text{subcycles}} \times \text{timestep} = 1\,\text{ps}$) is too small to promote equilibration by phase space exploration (Fig. 5.12B). Most of the paths in these ensembles consist of only a single phase point. Algorithmically, this is equivalent to paths containing three phase points, where the first and last points are not a part of the path ensemble (they belong to state A or B) and therefore excluded from being viable shooting points. This means that, for these 1-phase point paths, the configuration remains identical as only the momenta are modified. Thus, even though that many paths were accepted in these ensembles (Fig. 5.12A), almost no exploration of configurational phase space was performed.

The $P_A(\lambda|\lambda_A)$ discontinuities for $\lambda \geq 4.5\,\text{Å}$ values are due to the low amount of accepted paths ($\leq 100$) in the corresponding ensembles (Fig. 5.12A), even though approximately $500\,\text{ns}$ or more of shooting time was simulated in each of these ensembles (Fig. 5.12C).

It becomes clear that metastable states are present in the WT ABL-imatinib system, even for OP values as low as $\lambda = 2\,\text{Å}$, which introduces a separation of timescales that is *not* resolved by the one-dimensional $\lambda$ parameter. In other words, a return to state $A$ from
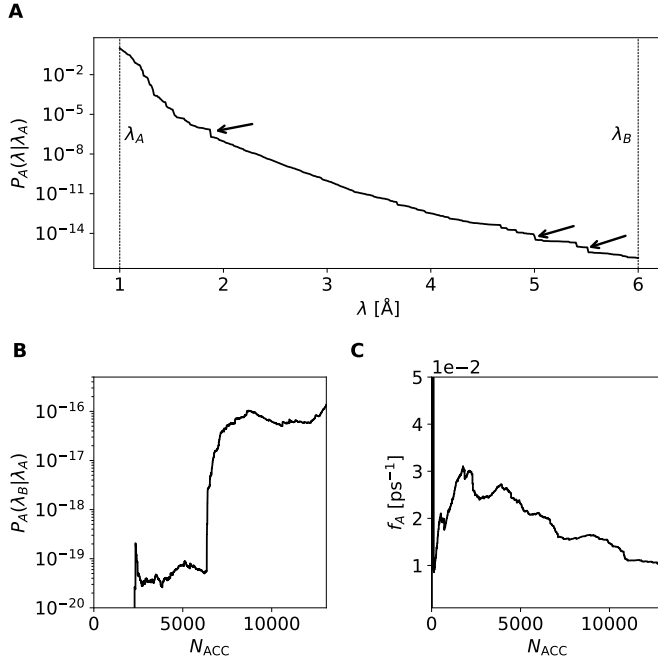
**Figure 5.11:** WT ABL ∞RETIS results. The first 5000 accepted paths were discarded to avoid initialization effects. **A**: The crossing probability profile $P_A(\lambda|\lambda_A)$. Arrows indicate visible discontinuities. **B**: The running estimate of the global crossing probability $P_A(\lambda_B|\lambda_A)$. **C**: The running estimate of the flux $f_A$. $N_{ACC}$ denotes the amount of accepted paths after removal of the first 5000.

$\lambda \geq 2\,\text{Å}$ can happen in a few (or even one) RETIS timesteps of 1 ps when no metastable state is encountered, or up to many thousands of timesteps when the path gets stuck in a metastable state.

**Figure 5.12:** Sampling analysis of the $\infty$RETIS WT ABL simulation. Only *shooting* cycles are included in the analysis, i.e. paths generated by the zero swap move in $[0^-]$ and $[0^+]$ are excluded. **A**: Amount of shooting moves performed $N_{\mathrm{sh}}(i)$ (black dots) and accepted $N_{\mathrm{acc}}(i)$ (black crosses) in each ensemble $E_i$. The acceptance ratio is shown in red. **B**: Average path lengths $\langle\tau\rangle_i$ of the ensembles $E_i$. Extremal phase points of paths are excluded. **C**: The total amount of simulated (shooting) time $\tau_{\mathrm{sh}}(i)$ (dots) in each ensemble $E_i$, and the accepted portion $\tau_{\mathrm{sh}}^{\mathrm{acc}}(i)$ thereof (crosses).

5.4 Conclusions

Both the RETIS and REPPTIS methodologies were not successful in sampling the dissociation pathways of imatinib from the ABL kinase domain. The main reason is rooted in the unbinding mechanism being more complex than the one-dimensional order parameter $\lambda$ can capture. The manner in which the resulting separation of timescales hinders the REPPTIS and RETIS methodologies, differs, and is now discussed.

For REPPTIS, the energy barriers orthogonal to $\lambda$ don't allow paths to equilibrate to the relevant regions of phase space. If the initial path is not truly representative of a dissociating pathway, the locality of the path ensemble definitions can hinder path sampling. This was observed especially for the ensembles close to the binding pocket, where the replica exchange mechanism to improve phase space exploration was rarely performed due to the dominance of LML paths. While RETIS significantly improved sampling for the $\lambda \in [2, 6]$ Å region, the separation of timescales would require a time step that is considerably smaller than 1 ps to allow sampling of the $\lambda \in [1, 2]$ Å region. This would require large disk spaces as paths related to larger $\lambda$ values were seen to exceed 60 ns, and longer overall simulation times as the MD simulation rate is reduced (by increased write frequency). There is, however, a larger problem regarding data interpretation, as $\lambda_B = 6$ Å is unlikely to represent *the* interface separating *the* long-lived metastable state between the deep binding pocket and the solvent.

Long MD simulations on the special purpose Anton 2 supercomputer have revealed that there are a multitude of long-lived metastable states, even close to and within the deep binding pocket (Roux et. al. [198], Shaw et. al. [210]). While the presence of metastable states is not problematic for REPPTIS, kinetic analysis of the ABL-imatinib system in the recent works of Refs [208, 209] revealed that the dissociation process is not well described by a one-dimensional reaction coordinate. The milestoning method was used by Elber et. al. [208], where Voronoi tessellation revealed a milestoning network with average connectivity of 2.93, whereas a value of 2 is expected for a one-dimensional description. Infrequent metadynamics simulations by Shekhar et. al. [209] used a 5-dimensional RC model to extract the rate constant of the ABL-imatinib system. The results of these studies were also included in Table 5.3, and are in much closer agreement with the few experimental values available.

Theoretically, the MC approach of generating paths is a strength of REPPTIS when orthogonal barriers are present. This is because consecutive path generation will see an unrepresentative initial path evolve into more representative ones (i.e. the importance sampling procedure), after which removal of the first thousands of paths in the analysis procedure can avoid initialization effects. In practice, however, this initial 'path equilibration' process can be hindered if the orthogonal barriers are too high. While a better initial path generation could have resulted in better convergence of the ABL-imatinib REPPTIS simulations, it is likely that the sampled paths would have only considered a fraction of the vast configurational network of metastable states. As such, the *accuracy* may still have been low, even if the *precision* was improved.

As the rugged energy landscape of many biological systems is expected to be similar to that of the ABL-imatinib systems considered in this chapter, there is a clear need for the REPPTIS methodology to better handle the presence of many metastable states, their (hidden) energy barriers, and their associated (hidden) timescale separations. A Voronoi tessellation approach could be introduced, where the resulting multidimensional PPTIS paths would be similar to those of the directional milestoning approach.

# 6

# EXTENDED REPPTIS METHODOLOGY

This chapter introduces a new methodology, replica exchange extended PPTIS (REPPEXTIS), that aims to enhance path memory and improve information exchange between path ensembles. As such, this chapter tackles **Research Objectives 1a** and **1b**, addressing the ergodic sampling issues and the accuracy limitations associated to limited path memory of PPTIS paths. The methodology is based on the MSM framework of Chapter 4, where REPPEXTIS paths inherently represent chains of connected PPTIS path segments. This path extension formalism increases path memory (longer paths), and simultaneously allows for paths to be (infinitely) swapped between REPPEXTIS ensembles.

The chapter starts with a short introduction, after which the new path ensembles are defined, the shoot-and-extend move is introduced, for which it is shown that detailed balance is maintained. Three strategies to perform path extensions are then discussed, after which the infinite swapping formalism and a multiple-replica formalism are detailed. The methodology is then applied to a shuffleboard potential to investigate memory enhancement, and to a two-dimensional rugged energy landscape to test the multiple-replica implementation. While the results are promising with regard to the memory enhancement, future work is required to make claims about improved convergence rates. This is discussed in the final section of the chapter, where

directions for future work are outlined, and a proposal for a multidimensional extension of the methodology is made.

## 6.1 INTRODUCTION

In the MSM framework of PPTIS, the elements of the transition matrix $M$ carried six indices

$$M_{ikl,i'k'l'} = P(S_i^{k,l} \to S_{i'}^{k',l'}), \qquad (6.1)$$

denoting the probability that a path of $\left[i^{\pm}\right]$ starting at $\lambda_{i+k}$ and ending at $\lambda_{i+l}$ will transition into a path of $\left[i'^{\pm}\right]$ starting at $\lambda_{i'+k'}$ and ending at $\lambda_{i'+l'}$ when propagated forward in time. As a PPTIS path can only extend into its neighboring ensembles, only $i' = i \pm 1$ entries of $M$ are non-zero. Moreover, only a few possibilities of $k'$ and $l'$ result in non-zero entires. An $\text{LMR}_{\left[i^{\pm}\right]}$ path, for example, can only transition into al $\text{LMR}_{\left[i+1^{\pm}\right]}$ or $\text{LML}_{\left[(i+1)^{\pm}\right]}$ path when propagated forwards in time. The replica exchange move in REPPTIS extends two paths of neighboring ensembles, where essentially extra path memory is gathered. This additional path memory could be used to extend the transition matrix $M$ to include 8 indices $M_{iklm,i'k'l'm'}^{\text{ext}}$, provided that sufficient swap moves are performed.

In principle, all the paths generated in a REPPTIS simulation could be extended in post-analysis, such that sufficient data is obtained to investigate the additional memory using the extended transition matrix. This is computationally expensive, and it would be desirable to have a methodology that actively uses such path extensions during the simulation. Such a methodology is presented here, where path extensions are coupled with an infinite swapping formalism. The methodology offers a trade-off between added memory and enhanced path space exploration.

## 6.2 THEORY

### 6.2.1 Path ensembles

An extended PPTIS (REPPEXTIS) $[i \pm N_{\text{ext}}]$ ensemble contains all *continuous* paths that can be decomposed into a chain of $N_{\text{ext}} + 1$ connected PPTIS path segments, where *at least one* of the segments is contained in the $\left[i^{\pm}\right]$ ensemble. Five example paths of the $[i \pm 3]$ ensemble are shown in Fig. 6.1, where it is seen that $[i \pm 3]$ paths span the $\lambda \in [\lambda_{i-4}, \lambda_{i+4}]$ region.
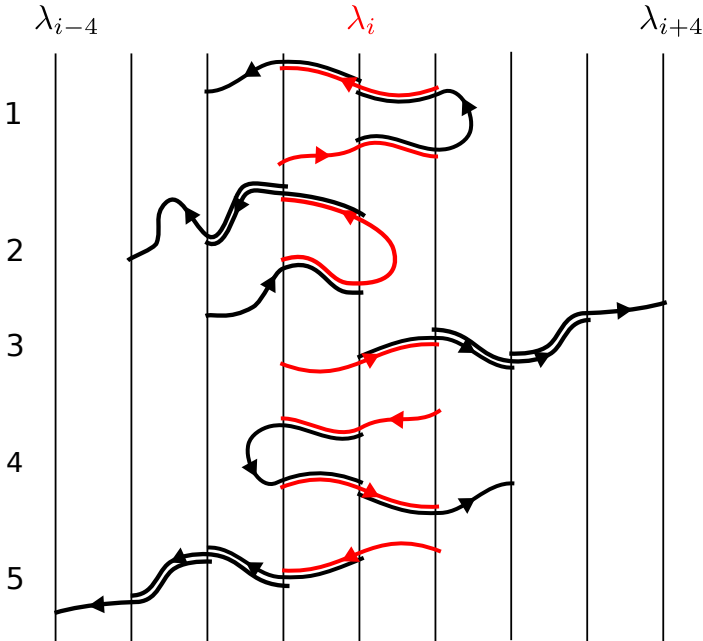
**Figure 6.1:** Five example paths of the $[i\pm3]$ ensemble. The paths are continuous, but can be viewed as $3+1 = 4$ PPTIS path segments, of which at least one is contained in the $[i^{\pm}]$ ensemble. Paths 1 and 4 each contain two $[i^{\pm}]$ segments, while paths 2, 3, and 5 contain only one $[i^{\pm}]$ segment. Path 3 is an example of a path that maximally extends to the right, while path 5 maximally extends to the left. The $\lambda$ region that $[i\pm3]$ paths can explore is therefore $[\lambda_{i-4}, \lambda_{i+4}]$.

Rather than listing all possible path types of a $[i\pm N_{\text{ext}}]$ ensemble, a constructive definition using the MSM interpretation of Chapter 4 is more instructive, which aligns well with the adapted 'shoot-and-extend' move that will be used to sample the $[i\pm N_{\text{ext}}]$ ensemble. Starting from a state $S_{\alpha}$ in the $[i^{\pm}]$ ensemble, the path is extended forwards or backwards in time until a new state $S_{\beta}$ is reached. This extension is similar as the one performed for the REPPTIS swapping move. If the propagation was forwards in time, the resulting segment chain is $S_{\alpha\beta} = (S_{\alpha}, S_{\beta})$. If the propagation was backwards, the segment chain is $S_{\beta\alpha} = (S_{\beta}, S_{\alpha})$. This chain is then again extended forwards or backwards in time until another state $S_{\gamma}$ is reached. This process is performed $N_{\text{ext}}$ times, resulting in the segment chain $S_{\alpha\prime\beta\prime\gamma\prime...}$, where the greek indices are reorganized according to the chosen time-directions. This chain of PPTIS states is now included in *one* specific REPPEXTIS state. Thus, using the MSM interpretation, the $[i\pm N_{\text{ext}}]$ ensemble contains all continuous paths that are a

chain $S_{\alpha\beta...}$ of $N_{\text{ext}}+1$ PPTIS path segments, where at least one index $n \in [\alpha, \beta, \ldots]$ corresponds to a state $S_n$ of the $[i^\pm]$ ensemble. This connection with the MSM interpretation allows for a straightforward analysis of REPPEXTIS simulations, where the transition matrix $M$ now considers states from one segment chain to another. The analysis is discussed further in Sec. 6.5.

As $[i\pm N_{\text{ext}}]$ ensembles were defined using $[i^\pm]$ ensembles, a REPPEXTIS simulation includes the same amount of path ensembles as a (RE)PPTIS simulation (for the same interfaces). However, two additional ensembles are added. These are the $[(N-1)'\pm N_{\text{ext}}]$ and $[N^-\pm N_{\text{ext}}]$ ensembles, which are state $B$ implementations of the $[0'\pm N_{\text{ext}}]$ and $[0^-\pm N_{\text{ext}}]$ ensembles. In its current state, the methodology is therefore only applicable to bounded $A \rightleftharpoons B$ reactions. The addition of the $[(N-1)'\pm N_{\text{ext}}]$ and $[N^-\pm N_{\text{ext}}]$ ensembles is necessary to ensure that paths can still be extended when they reach the $\lambda_B$ interface.

### 6.2.2 Shoot-and-extend move

Consider the REPPEXTIS path $o = \{o_3, o_2, o_1, o_4\}$ of the $[i\pm 3]$ ensemble at the top of Fig. 6.2A. The path segments $o_i$ are ordered by forwards progression in time, where the path is seen to be a chain of segments going through $\left\{ [i^\pm], [(i+1)^\pm], [i^\pm], [(i-1)^\pm] \right\}$.

The shooting move in $[i\pm 3]$ consists of shooting from a $[i^\pm]$ path segment. The first step is thus to select a path segment from the $[i\pm 3]$ path, after which a regular PPTIS shooting move in $[i^\pm]$ is performed. For the path of Fig. 6.2, segment $o_3$ is selected, while $o_1$ would also be a viable choice. The newly created $[i^\pm]$ path segment is then extended by propagating the segment $N_{\text{ext}}$ times in randomly chosen time directions. To allow for an infinite swapping formalism, these propagation directions should not be made deterministically. If, for example, one were to only propagate forwards in time, then the $[i\pm 3]$ path ensemble would be restricted to path chains that start with an $[i^\pm]$ segment. Such path chains are then, by construction, only a part of the $[i\pm 3]$ ensemble. A time-reversal move could be combined with a swap attempt to allow for swapping without MD integration, but it is better to aim for less restricted replica exchange. The newly generated path in Fig. 6.2 was propagated twice forwards in time, and once backwards in time, resulting in the new $[i\pm 3]$ path $n = \{n_1, n_2, n_3, n_4\}$.

To obey detailed balance, Eq. 2.20 must hold, and if the Metropolis-Hastings acceptance scheme is used, the acceptance

**Figure 6.2:** Shoot-and-extend move from a $[i\pm3]$ path $o$ to a new $[i\pm3]$ path $n$ and the reverse operation.

probability is given by Eq. 2.21. The probability $P_{\mathrm{acc}}^{o\to n}$ of generating a (n)ew $[i\pm N_{\mathrm{ext}}]$ path from an (o)ld path is now discussed, where an example showing both the $o$ to $n$ and the $n$ to $o$ transitions is shown in Fig. 6.2. To not overload the notation, a $[i\pm N_{\mathrm{ext}}]$ superscript or subscript is not included in the following equations.

The probability of generating a new path $n$ from an old path $o$ is given by

$$
\begin{aligned}
P_{\text{gen}}^{o \to n} = 1 & \\
\overset{(1)}{\times}\ & P_{\text{sel, seg}}^{o \to o_i} \\
\overset{(2)}{\times}\ & P_{\text{sel, ph}}^{o_i \to o_i(k)} \\
\overset{(3)}{\times}\ & P_{\text{mod, ph}}^{o_i(k) \to n_j(l)} \\
\overset{(4)}{\times}\ & P_{\text{engine, fw+bw}}^{n_j(l) \to n_j} \\
\overset{(5)}{\times}\ & P_{\text{sel, direcs } [d_1^n,...,d_{N_{\text{ext}}}^n]}^{n_j \xrightarrow{\text{allow}} n} \\
\overset{(6)}{\times}\ & \prod_{q=1}^{N_{\text{ext}}} P_{\text{engine}, d_q^n}^{\{n_s\}_{s=j}^{j+q-1} \to \{n_s\}_{s=j}^{j+q}},
\end{aligned}
\tag{6.2}
$$

where the notation of each step is explained in detail below. The probability of generating the reverse transition is

$$
\begin{aligned}
P_{\text{gen}}^{n \to o} = 1 & \\
\overset{(1)}{\times}\ & P_{\text{sel,seg}}^{n \to n_j} \\
\overset{(2)}{\times}\ & P_{\text{sel, ph}}^{n_j \to n_j(l)} \\
\overset{(3)}{\times}\ & P_{\text{mod, ph}}^{n_j(l) \to o_i(k)} \\
\overset{(4)}{\times}\ & P_{\text{engine, fw+bw}}^{o_i(k) \to o_i} \\
\overset{(5)}{\times}\ & P_{\text{sel, direcs } [d_1^o,...,d_{N_{\text{ext}}}^o]}^{o_i \xrightarrow{\text{allow}} o} \\
\overset{(6)}{\times}\ & \prod_{q=1}^{N_{\text{ext}}} P_{\text{engine}, d_q^o}^{\{o_s\}_{s=i}^{i+q-1} \to \{o_s\}_{s=i}^{i+q}}.
\end{aligned}
\tag{6.3}
$$

We now cover the $o \to n$ generation in detail. The first step is to select a path segment of the $[i^{\pm}]$ ensemble from the $o$ path. The probability of choosing specifically path segment $o_i$ from the $o$ path is denoted by $P_{\text{sel, seg}}^{o \to o_i}$, resulting in the first (1) probability multiplier of Eq. 6.2. Different strategies exist to select path segments, and the notation is kept general for now. Steps (2) to (4) are identical to a PPTIS shooting move, where selection of phase point $o_i(k)$ from path segment $o_i$ is denoted by $P_{\text{sel, ph}}^{o_i \to o_i(k)}$, the modification of phase

point $o_i(k)$ to $n_j(l)$ is denoted by $P_{\text{mod, ph}}^{o_i(k)\rightarrow n_j(l)}$, and the forward (fw) and backwards (bw) engine propagation to create a new $\left[i^{\pm}\right]$ path $n_j$ from the modified phase point $n_j(l)$ is denoted by $P_{\text{engine, fw+bw}}^{n_j(l)\rightarrow n_j}$. The next steps (5) and (6) concern the path extensions. The selection of time directions (5) and the actual engine propagations (6) are separated for generality. Depending on the strategy, the probability of selecting $N_{\text{ext}}$ time directions $[d_1^n, \ldots, d_{N_{\text{ext}}}^n]$ (which are either fw or bw) can be different. The probability of selecting a set of time directions that *allows* the $n_j$ path segment to be extended into the $n$ path is given by $P_{\text{sel, direcs } [d_1^n,...,d_{N_{\text{ext}}}^n]}^{n_j \xrightarrow{\text{allow}} n}$. Selecting 3 backwards extensions in the example of Fig. 6.2 can for instance never result in the $n$ path (as it requires 2 forwards and 1 backwards extension). A shorthand notation will be used later on to alleviate the notation, where the selection of allowed time directions will be denoted as $P_{\text{sel, direcs}}^{n_j \xrightarrow{\text{allow}} n}$. The *allowed* part was emphasized, because the actual path propagations can still deviate from creating the $n$ path if a stochastic engine is used. The actual time extensions are performed in the final step (6), where the total probability of these extensions is given by the product of the individual path segment propagations. The notation $P_{\text{engine},d_q^n}^{\{n_s\}_{s=j}^{j+q-1}\rightarrow\{n_s\}_{s=j}^{j+q}}$ denotes the probability of extending the (continuous) path of segments $\{n_s\}_{s=j}^{j+q-1}$ into the extended path of segments $\{n_s\}_{s=j}^{j+q}$ by propagation in time direction $d_q^n$. It is implied here that a backwards time propagation results in the segment being appended at position 0 of the path chain, while forward propagation results in the segment being appended at final position. In short, this will be denoted as $P_{\text{engine},d_q^n}^{\text{ext}}$, being the propagation probability of the $q$th segment extension in the $n$ path. The generation probability of going from the new path $n$ to the old path $o$ is given by Eq. 6.3, where the same steps are taken.

Plugging these equations into the Metropolis-Hastings acceptance

criterion (Eq. 2.21), we get

$$
P_{\text{acc}}^{o \to n} = \min \left\{ 1, \frac{\mathcal{P}(n) P_{\text{gen}}^{n \to o}}{\mathcal{P}(o) P_{\text{gen}}^{o \to n}} \right\}
$$

$$
= \min \left\{ 1, \frac{P_{\text{sel, ph}}^{n_j \to n_j(l)}}{P_{\text{sel, ph}}^{o_i \to o_i(k)}} \right.
$$

$$
\times \frac{\mathcal{P}(n) \times P_{\text{mod, ph}}^{n_j(l) \to o_i(k)} \times P_{\text{engine, fw+bw}}^{o_i(k) \to o_i} \times \prod_{q=1}^{N_{\text{ext}}} P_{\text{engine},d_q^o}}{\mathcal{P}(o) \times P_{\text{mod, ph}}^{o_i(k) \to n_j(l)} \times P_{\text{engine, fw+bw}}^{n_j(l) \to n_j} \times \prod_{q=1}^{N_{\text{ext}}} P_{\text{engine},d_q^n}}
$$

$$
\left. \times \frac{P_{\text{sel, seg}}^{n \to n_j} \times P_{\text{sel, direcs}}^{o_i \to o}}{P_{\text{sel, seg}}^{o \to o_i} \times P_{\text{sel, direcs}}^{n_j \to n}} \right\},
$$

$$\tag{6.4}$$

where the fraction was rearranged and split into a product of three fractions. The middle fraction cancels out completely. Assuming microscopic reversibility (Eq. 2.8), it is possible to define the probability of a path $k$ as $\mathcal{P}(k) = \rho(k^{\text{ph}}) P_{\text{int}}(k | k^{\text{ph}})$, where $P_{\text{int}}(k | k^{\text{ph}})$ denotes the probability of creating path $k$ by going backwards and forwards in time from phasepoint $k^{\text{ph}}$. This is true for any general path, not just PPTIS or TIS paths. It is assumed that, if ph is the $i$-th phasepoint in a path of $N$ phasepoints, then $i - 1$ backwards steps and $N - i$ forwards steps are to be taken. Since the product of all engine probabilities concerning $o$ path segments create one continuous path (the overlapping parts are not integrated twice, where path extensions are performed similarly to a one part of REPPTIS swap move), their product is simply $P_{\text{int}}(o | o_i(k)) = \mathcal{P}(o) / \rho(o_i(k))$. This cancels the appearance of $\mathcal{P}(o)$ in the denominator of Eq. 6.4. In similar fashion, the $\mathcal{P}(n)$ term in the nominator is cancelled. The middle fraction is then reduced to

$$
\frac{P_{\text{mod, ph}}^{n_j(l) \to o_i(k)} \rho(n_j(l))}{P_{\text{mod, ph}}^{o_i(k) \to n_j(l)} \rho(o_i(k))}.
$$

Assuming phasepoint modification is symmetric, the phasepoint modification probabilities cancel out. For aimless shooting, the kinetic contributions to the phasepoint energies is equal, and since their configurational energies are identical (*if* by construction only the velocities are modified), the phasepoint densities $\rho$ also cancel out.

As such, the middle fraction cancels out, and we have that

$$P_{\text{acc}}^{o \to n} = \min \left\{ 1, \frac{P_{\text{sel, ph}}^{n_j \to n_j(l)}}{P_{\text{sel, ph}}^{o_i \to o_i(k)}} \times \frac{P_{\text{sel, seg}}^{n \to n_j} \times P_{\text{sel, direcs}}^{o_i \to o}}{P_{\text{sel, seg}}^{o \to o_i} \times P_{\text{sel, direcs}}^{n_j \to n}} \right\}. \qquad (6.5)$$

We will equally well obey detailed balance [113] when this acceptance probability is split in a product

$$P_{\text{acc}}^{o \to n} = \min \left\{ 1, \frac{P_{\text{sel, ph}}^{n_j \to n_j(l)}}{P_{\text{sel, ph}}^{o_i \to o_i(k)}} \right\} \times \min \left\{ 1, \frac{P_{\text{sel, seg}}^{n \to n_j}}{P_{\text{sel, seg}}^{o \to o_i}} \right\} \times \min \left\{ 1, \frac{P_{\text{sel, direcs}}^{o_i \to o}}{P_{\text{sel, direcs}}^{n_j \to n}} \right\},$$
$$(6.6)$$

allowing early rejection strategies. Three different strategies for segment selection and extension are now discussed, rising in complexity. The first strategy is a naive approach, where any segment of the $o$ path is selected with equal probability. The second and third strategies enforce selection of $[i^{\pm}]$ segments only. The second strategy chooses random time directions, while the third strategy chooses time directions that favor equal amounts of forwards and backwards extensions.

### 6.2.3 First strategy

A first strategy is to select any of the segments of the $o$ path with equal probability $1/N_{\text{ext}}$. If the selected segment is not a part of the $[i^{\pm}]$ ensemble, the move is rejected. The segment selection probability thus becomes

$$P_{\text{sel, seg}}^{o \to o_i} = \theta_{[i^{\pm}]}(o_i) \frac{1}{N_{\text{ext}} + 1}, \qquad (6.7)$$

where a factor $\theta_{[i^{\pm}]}(o_i)$, the indicator function that is 1 if the selected segment is a part of the $[i^{\pm}]$ ensemble, and 0 otherwise, is to be added to Eq. 6.5. The selection of time directions is then chosen uniformly, where the probability of doing $q \in \{1, \ldots, N_{\text{ext}}\}$ forward extensions all have probability $1/N_{\text{ext}}$. The direction selection probability is then

$$P_{\text{sel, direcs}}^{o_i \to o} = \frac{1}{N_{\text{ext}}}. \qquad (6.8)$$

The phasepoint selection probabilities are given by $\frac{1}{N^{n_j}}$ and $\frac{1}{N^{o_i}}$, where $N^{n_j}$ and $N^{o_i}$ are the amount of phase points in the $n_j$ and $o_i$ path segments, respectively. The resulting acceptance probability of Eq. 6.5 then simply becomes

$$P_{\text{acc}}^{o \to n} = \theta_{[i^{\pm}]}(o_i) \min \left\{ 1, \frac{N^{o_i}}{N^{n_j}} \right\}. \qquad (6.9)$$

Which is exactly the same as the regular $[i^\pm]$ shooting move acceptance criterion, where now an additional check is performed for acceptable segment selection. In regions of steep energy barriers, LML or RMR paths will be dominant (depending on which side of the barrier the ensemble is located), where these paths will likely continue to 'roll down' this barrier. As such, the extensions are likely not be a part of the $[i^\pm]$ ensemble, and most shooting moves will be automatically rejected by the $\theta_{[i^\pm]}$ criterion. Such ensembles were seen to be troublesome for the biological applications, and it is therefore imperative to sufficiently sample these regions. This can be fixed (a) by adopting an asynchronous sampling strategy [118, 119], where these quickly-rejected ensembles will be sampled more frequently, and/or (b) by only considering $[i^\pm]$ segments for selection. which is done in the second and third strategy.

### 6.2.4 Second strategy

Considering only $[i^\pm]$ segments for selection, the segment selection probability becomes path dependent. We first discuss the simplest selection rule, before covering the third strategy used in the simulations below. Denote $N^o_{[i^\pm]}$ and $N^n_{[i^\pm]}$ as the amount of $[i^\pm]$ segments in the $o$ and $n$ paths, respectively. The simplest selection rule is to select either of the segments with equal probability, where

$$
\begin{aligned}
P^{o \to o_i}_{\text{sel, seg}} &= \frac{1}{N^o_{[i^\pm]}}, \\
P^{n \to n_j}_{\text{sel, seg}} &= \frac{1}{N^n_{[i^\pm]}}.
\end{aligned}
\tag{6.10}
$$

Keeping the selection of time directions probabilities uniform, the acceptance probability becomes

$$
P^{o \to n}_{\text{acc}} = \min \left\{ 1, \frac{N^{o_i}}{N^{n_j}} \frac{N^o_{[i^\pm]}}{N^o_{[i^\pm]}} \right\}.
\tag{6.11}
$$

As $N^n_{[i^\pm]}$ cannot be known beforehand, the extensions have to be performed before the acceptance criterion can be evaluated. It is, however, possible to use alternative path ensemble definitions for which the $N_{[i^\pm]}$ terms disappear from the acceptance criterion. This strategy maximizes the acceptance probability, and has been used in path sampling moves involving subtrajectories, like the stone skipping, web throwing, and wire fencing moves [129, 130]. The idea

is discussed below, while an in-depth explanation can be found in the supplementary information of Ref. [129]. The alternative path ensemble uses a path probability $\tilde{\mathcal{P}}$

$$\tilde{\mathcal{P}}(k) = w(k)\mathcal{P}(k), \tag{6.12}$$

where a path $k$ is now weighted by $w(k)$.

Assume that we simply remove the $N^o_{[i\pm]}$ and $N^n_{[i\pm]}$ term from the acceptance criterion in Eq. 6.11, then the original path ensemble is no longer sampled. Instead, one would would sample the alternative path ensembles with weights $w_{[i\pm N_{\text{ext}}]}(k) = N^k_{[i\pm]}$. To see this, the new path probabilities $\tilde{\mathcal{P}}$ are used in the Metropolis acceptance criterion (Eq. 2.21), where now

$$P^{o\rightarrow n}_{\text{acc}} = \min\left\{1, \frac{\tilde{\mathcal{P}}(n)P^{n\rightarrow o}_{\text{gen}}}{\tilde{\mathcal{P}}(o)P^{o\rightarrow n}_{\text{gen}}}\right\}$$

$$= \min\left\{1, \frac{\mathcal{P}(n)\left(N^n_{[i\pm]}\right)P^{n\rightarrow o}_{\text{gen}}}{\mathcal{P}(o)\left(N^o_{[i\pm]}\right)P^{o\rightarrow n}_{\text{gen}}}\right\}. \tag{6.13}$$

The only change is the presence of these weights, which cancel the $N^o_{[i\pm]}$ and $N^n_{[i\pm]}$ stemming from the generation probabilities in Eq. 6.11. As such, the acceptance probability again becomes equal to the regular PPTIS shooting move acceptance criterion. The subscript $[i\pm N_{\text{ext}}]$ was explicitly added to denote that the weight of a path $k$ now depends on the $[i\pm N_{\text{ext}}]$ ensemble. This has to be taken into account when swapping paths between the ensembles, as discussed in the next section.

In post analysis, the path probabilities have to corrected by the inverse weights to obtain correct statistics. For the simulation with the adapted path ensemble, the paths are sampled with probability $\tilde{\mathcal{P}}$, where the original ensemble average $\langle A \rangle$ of a property $A$ is recovered by

$$\langle A \rangle = \frac{\sum_{j=1}^{N_{\text{paths}}} A(j)\, w^{-1}(j)}{\sum_{j=1}^{N_{\text{paths}}} w^{-1}(j)}, \tag{6.14}$$

where $A(j)$ is the property evaluated in path $j$, and $N_{\text{paths}}$ is the amount of sampled paths.

### 6.2.5   Third strategy

The idea of the third strategy is to achieve a balance between forward and backward extensions. This 'temporal' balance correlates with 'spatial' balance, as disfavoring extension in one time direction only decreases the likelihood of paths extending all the way to $\lambda_{i\pm(N_{\text{ext}}+1)}$. As such, this strategy produces $[i\pm N_{\text{ext}}]$ paths that focus exploration around the $\lambda_i$ interface. Depending on the underlying dynamics, this strategy can be favorable. Consider for example a steep energy barrier region, where both forward and backwards extensions of a $\left[i^{\pm}\right]$ path segment are most likely to roll down the barrier. Long extensions in one time direction are then unlikely to result in swaps, as the segments of the ensembles 'down the hill' are unlikely to climb all the way back up to $\lambda_i$. As such, it may be favorable to keep paths around the $\lambda_i$ interface, where the chance of being swapped with direct neighboring ensembles is higher.

To achieve this, the third strategy uses a binomial distribution to select the amount of forward time extensions

$$P_{\text{sel, direcs}}^{o_i \to o} = \binom{N_{\text{fw}}^{o_i \to o}}{N_{\text{ext}}}, \tag{6.15}$$

where $N_{\text{fw}}^{o_i \to o}$ is the amount of forward extensions that are required for (allowed) creation of the $o$ path from the $o_i$ segment.

This is, however, a *segment dependent property* rather than an $[i\pm N_{\text{ext}}]$ path property. As such, the segment selection probability has to be changed such that a *path dependent property* arises and the *path*-reweighting trick can be applied for high acceptance. To this end, the segment selection probability is defined as

$$P_{\text{sel, seg}}^{o \to o_i} = \frac{\binom{N_{\text{fw}}^{o_i \to o}}{N_{\text{ext}}}}{\sum_{r=1}^{N_{\left[i^{\pm}\right]}^{o}} \binom{N_{\text{fw}}^{o_r \to o}}{N_{\text{ext}}}}, \tag{6.16}$$

where segments $o_i$ that have a larger probability of being 'binomially extended' into path $o$ are more likely to be selected. Plugging the extension and selection probabilities into Eq. 6.5, the acceptance

probability becomes

$$P_{\text{acc}}^{o \to n} = \min \left\{ 1, \frac{N^{o_i}}{N^{n_j}} \right.$$

$$\times \frac{\dfrac{\binom{N_{\text{fw}}^{n_j \to n}}{N_{\text{ext}}}}{\sum_{r=1}^{N_{[i\pm]}^n} \binom{N_{\text{fw}}^{n_r \to n}}{N_{\text{ext}}}}}{\dfrac{\binom{N_{\text{fw}}^{o_i \to o}}{N_{\text{ext}}}}{\sum_{r=1}^{N_{[i\pm]}^o} \binom{N_{\text{fw}}^{o_r \to o}}{N_{\text{ext}}}}} \tag{6.17}$$

$$\left. \times \frac{\binom{N_{\text{fw}}^{n_j \to n}}{N_{\text{ext}}}}{\binom{N_{\text{fw}}^{o_i \to o}}{N_{\text{ext}}}} \right\},$$

where the binomial coefficients of the direction selection cancel with the nominators of the segment selection probabilities to give

$$P_{\text{acc}}^{o \to n} = \min \left\{ 1, \frac{N^{o_i}}{N^{n_j}} \times \frac{\sum_{r=1}^{N_{[i\pm]}^o} \binom{N_{\text{fw}}^{o_r \to o}}{N_{\text{ext}}}}{\sum_{r=1}^{N_{[i\pm]}^n} \binom{N_{\text{fw}}^{n_r \to n}}{N_{\text{ext}}}} \right\}. \tag{6.18}$$

A path $k$ is now assigned the weight

$$w_{[i\pm N_{\text{ext}}]}(k) = \sum_{r=1}^{N_{[i\pm]}^k} \binom{N_{\text{fw}}^{k_r \to k}}{N_{\text{ext}}}, \tag{6.19}$$

which denotes the combined probability of creating this path from its $[i\pm]$ segments using binomially distributed time extensions. The subscript $[i\pm N_{\text{ext}}]$ was explicitly added to denote that the weight of a path $k$ now depends on the $[i\pm N_{\text{ext}}]$ ensemble. As mentioned before, this has to be taken into account when swapping paths between the $[i\pm N_{\text{ext}}]$ ensembles, which is now discussed.

## 6.3 Infinite swapping

The infinite swapping formalism was introduced in Refs. [118, 119], and can straightforwardly be applied to the $[i\pm N_{\text{ext}}]$ path ensembles. The key idea is repeated here, and the reader is referred to these works for more in-depth derivations and discussions. The asynchronous formalism was also introduced in these works, where its implementation is technical and not (yet) considered here.

An infinite swapping formalism for a REPPTIS simulation was impossible, as the swap move itself involved MD integration. In the novel path sampling methodology, the extension of paths is integrated within the shooting move, where swaps can now be performed freely. An $[i\pm N_{\text{ext}}]$ ensemble was defined to contain all paths that can be decomposed into a chain of $N_{\text{ext}}+1$ connected PPTIS path segments, where at least one of the segments is a part of the $[i^{\pm}]$ ensemble. A $[i\pm N_{\text{ext}}]$ path is therefore a part of all the $[j\pm N_{\text{ext}}]$ ensembles for which it contains at least one $[j^{\pm}]$ path segment. Depending on the shooting strategy, the weight of a path can be different for different $[j\pm N_{\text{ext}}]$ ensembles.

Shooting and (infinite) swapping moves are performed in alternating fashion. First, shooting moves are performed in all $[i\pm N_{\text{ext}}]$ ensembles. Then, an infinite swapping move is performed, after which the paths are redistributed over ensembles. This redistribution is done according to the probability matrix $P$ of path locations following an infinite number of swapping moves. The elements $P_{ij}$ of this matrix denote the probability that a path $k_i$ ends up in ensemble $[j\pm N_{\text{ext}}]$ after infinite swaps. The calculation of $P$ is done using a permenant representation

$$P = \frac{W_{ij}\text{perm}\big(W[ij]\big)}{\text{perm}(W)}, \tag{6.20}$$

where $W$ is the weight matrix with elements $W_{ij}$ denoting the weight of path $k_i$ in ensemble $[j\pm N_{\text{ext}}]$, $\text{perm}(W)$ is the permenant of $W$, and $\text{perm}\big(W[ij]\big)$ is the permenant of $W$ with the $i$-th row and $j$-th column removed.

## 6.4   Multiple replicas

As seen in the biological applications, the REPPTIS path ensembles can get stuck in regions of path space by bad initial path selection. More generally, the complex rugged free energy landscape of biological systems can be expected to have multiple reactive pathways, where *accurate* convergence can be slow due to paths being stuck in one or few pathways. As such, it could be beneficial to include $N_l$ multiple replicas (or 'levels') of the same path ensemble. If $N_l$ initial paths are generated, and $N_l$ replicas of each path ensemble are run concurrently, the sampling could potentially be superior to running $N_l$ independent simulations due to the efficient information exchange of the infinite swapping formalism. Additional replicas could be introduced specifically in difficult-to-sample $\lambda$-intervals. Furthermore,

if shoot-and-extend moves are performed in parallel (or asynchronously), it would allow for usage of more hardware.

A replica of an existing $[i\pm N_{\text{ext}}]$ ensemble follows the exact same definition as the original ensemble. As such, paths can therefore automatically be 'self-swapped' between replicas representing the same $[i\pm N_{\text{ext}}]$ ensemble. The added benefit of multiple replicas could lie in increased exchange of information between *different* $[i\pm N_{\text{ext}}]$ ensembles. As $N_l$ replicas increase the probability of observing a specific path type by a factor $N_l$, and as the infinite swapping formalism *will* exchange path information whenever it is possible, this hints towards a potential increase in sampling efficiency. However, as the probability of self-swapping also increases, this enhancement may be diminished. It is thus unclear whether the use of multiple replicas will introduce improvements other than (a) allowing for extra hardware in combination with parallel MD integrator schemes and (b) additional sampling of slow-converging ensembles by only introducing replicas for these ensembles. Investigating the potential benefits of multiple-replicas is an interesting avenue for future work.

## 6.5 ANALYSIS

As noted in the introduction, the MSM formalism of Chapter 4 can be adapted and applied to REPPEXTIS. The states $S_{\alpha\prime} = S_{\alpha\beta\gamma\dots}$ of the REPPEXTIS MSM are defined as all *possible* $N_{\text{ext}} + 1$ walks $(S_\alpha, S_\beta, S_\gamma, \dots)$ in the PPTIS MSM. For REPPEXTIS, state $S_A$ is defined as the set of all states $S_{\alpha\beta\dots\omega}$ whose final segment $S_\omega$ is the $S_{[0^-]}^{1,1}$ state, which corresponds to the (RMR) paths of the $[0^-]$ ensemble. The state $S_{\mathcal{B}}$ is defined as the set of all states $S_{\alpha\beta\dots\omega}$ whose final segment $S_\omega$ is the $S_{([N-2^\pm])}^{-1,1}$ state, which corresponds to the LMR paths of the $[(N-1)^\pm]$ ensemble. The transition probability between states $S_{\alpha\beta\dots\omega}$ and $S_{\beta\dots\omega\psi}$ is then defined by the fraction of $S_{\beta\dots\omega*}$ paths that have been observed to transition towards $S_{\beta\dots\omega\psi}$, where $*$ can be either of the two possible PPTIS path segments the $S_{\beta\dots\omega}$ chain can transition to (in the forwards time-direction). As such, the memory of the MSM states is defined by the first $N_{\text{ext}}$ segments of its $N_{\text{ext}}+1$-segment chain. As multiple REPPEXTIS ensembles sample the same paths, there are multiple ensembles estimators for one transition. For now, the transition probability is simply defined as the weighted average of the ensemble estimators. It could be interesting to investigate differences in the ensemble estimations for a single transition, but such analysis is not considered here. It could also be of interest to

truncate REPPEXTIS paths in post-analysis, simply by considering only the last $N' < N_{\text{ext}}$ segments of the paths as their memory. This could potentially serve as a tool to investigate memory effects.

The crossing probability profile $P_(\lambda_i|\lambda_A)$ is extracted using the same analysis procedure as defined in Chapter 4 (and detailed in **paper III**). To calculate the hitting probabilities related to the $\lambda_i$ interface, the absorbing state at $\lambda_i$ is defined as the set of all states $S_{\lambda_i} = S_{\alpha\beta\ldots\omega}$ whose final segment $S_\omega$ is the $S^{-1,1}_{[(i-1)^\pm]}$ state, which corresponds to the LMR paths of the $\left[(i-1)^\pm\right]$ ensemble. As state $B$ is now included in the ensembles, the stationary distribution $\pi$ can be calculated, which is required to properly weight the set of states considered in $S_A$.

## 6.6 IMPLEMENTATION

REPPEXTIS was implemented in Python, together with an implementation of RETIS and REPPTIS for comparison. The code is based on the PyRETIS structure, where the definitions for the ensemble objects, path objects, and engine propagation functions were redesigned with flexibility in mind. The code is available at GitHub at `https://github.com/WouterWV/reppextis`. The fast permanent calculation scripts were adapted from $\infty$RETIS [118, 119].

## 6.7 APPLICATIONS

The methodology is tested on two systems, designed to test the hypothesized benefits of the methodology: increased convergence rate due to increased replica exchange, and increased memory due to the longer paths. First, a shuffle-board potential is used to show the memory increase of the methodology, and second, a rough energy landscape with metastable states was used to test the inclusion of multiple replicas.

### 6.7.1 Shuffleboard potential

A shuffleboard potential was used to test the memory increase of REPPEXTIS. Such potential was previously used by van Erp to discuss the effects of RC-dependency of path sampling methods [133]. The shuffleboard potential (Fig. 6.3) differs from his set-up, as the barrier considered is cosine-shaped rather than a triangular energy barrier (i.e. linear increase to the peak from both sides).

A two-dimensional Langevin particle is confined to a rectangular box with dimensions $[L_x, L_y] = [0.9, 0.3]$, where reduced units are used (see below).
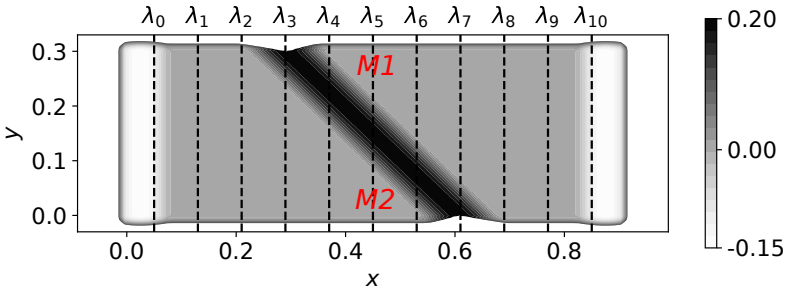
**Figure 6.3:** The shuffleboard potential. The barrier is a cosine-shaped bump of height $V_m = 0.2$ that extends along a $y = \theta x$ line, with $\theta = -1$. Regions $M1$ and $M2$ are causal to tunneling effects for low $N_{\text{ext}}$ values (see text).

The walls of the box are harmonic potentials ($k = 2000$) that are only felt when the particle goes beyond the box boundaries. The center of the harmonics are located at the wall edges ($x = 0, x = 0.9, y = 0, y = 0.3$), such that the transition is smooth. Cosine dips (bump) are used to model states $A$ and $B$ (barrier $M$). For the states $A$ and $B$, the potential is given by $V_S(x) = V_S \cos(2\pi(x - x_S)/L)$ if $|x - x_S| < L/4$, else 0. Parameters used are $L = 0.2$ and $V_S = -0.15$ for $S = A, B$, and $x_S = 0.025 (S = A)$ or $0.875 (S = B)$. The barrier $M$ goes through the midpoint of the box, and is angled at a slope $\theta = -1$. The potential is given by $V_M(x, y) = V_m \cos(2\pi D(x, y)/L_m)$ if $D(x, y) < L_m/4$ else 0, with $D(x, y)$ the distance of the particle to the line $y = \theta x + 0.6$. Parameters used are $V_m = 0.2$, $L_m = 0.2$ and $\theta = -1$. An additional overall slope $V_{\text{slope}}(x) = 0.2x$ is added, slightly pulling the particle towards state $A$.

The Langevin dynamics were integrated using reduced units in a Lennard-Jones type of unit system based on argon [211]. The parameters used are: temperature $T = 0.05$, particle mass $m = 0.1$, friction coefficient $\gamma = 25$, and integration timestep $\Delta t = 0.002$. Ten interfaces are positioned equidistantly, with $\lambda_A = 0.05$ and $\lambda_B = 0.85$. REPPEXTIS simulations were performed with $N_{\text{ext}} = 0, 1, 2, 3, 4, 5$, and 6 segment extensions, using the third strategy (binomial extension) for the shoot-and-extend moves. The $N_{\text{ext}} = 0$ simulation is essentially a PPTIS simulation. A RETIS simulation was performed as reference, where 75 % of the MC moves were shooting moves, and

the other 25 % swapping moves. All simulations ran concurrently for
~10 hours on a dell XPS 15 7590 laptop. Simulation cycle numbers
(in thousands) were (100, 100, 80, 80, 70, 60, 50) for the $N_{ext}$ val-
ues (0, 1, 2, 3, 4, 5, 6), respectively, where 1 cycle denotes 1 MC
move in all ensembles and 1 infinite swap. The RETIS simulation
was performed for 100 000 cycles.

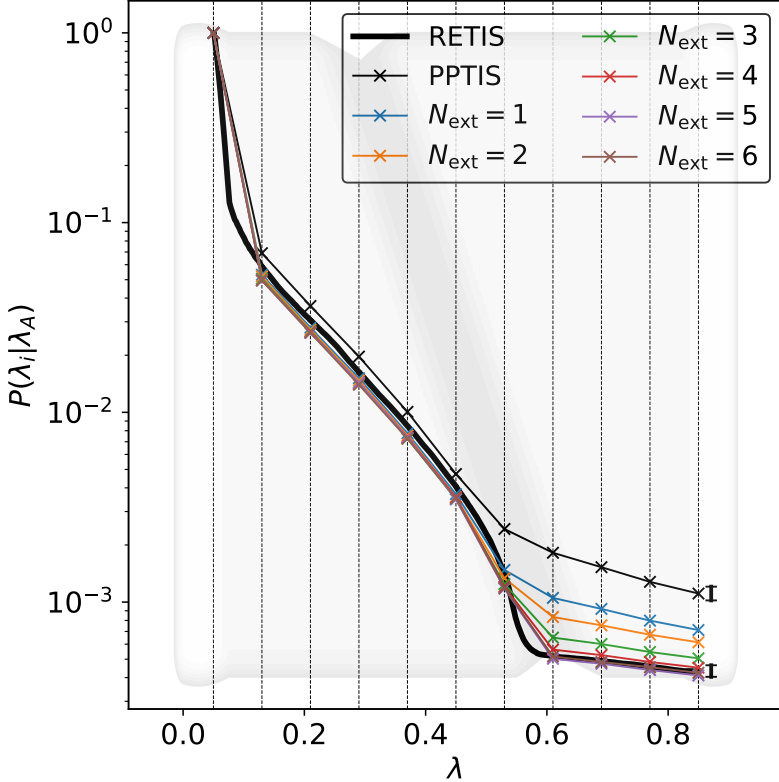The results are shown in Fig. 6.4. PPTIS has the largest overes-



**Figure 6.4:** The crossing probability profiles $P_A(\lambda_B|\lambda_A)$ for the RETIS,
PPTIS, and REPPEXTIS simulations. The shuffleboard potential is shown
as a shade. The two error bars correspond to the standard error of
$P_A(\lambda_B|\lambda_A)$ estimates by the PPTIS and RETIS simulations. An error es-
timate calculation for the REPPEXTIS simulations is not yet implemented.

timation of the global crossing probability, where so-called (apparent)
*tunneling effects* [121] are largest. For PPTIS paths in the $\left[5^{\pm}\right]$ en-
semble, the $p^{\pm}$ local crossing probability is largely overestimated due
to their sampling at the locations $M1$ and $M2$, respectively (Fig. 6.3).
At both of these locations, paths going from their left interface to

right interface (or vice-versa) do not require crossing of the energy barrier. Furthermore, this is true for every $\left[i^{\pm}\right]$ ensemble along the barrier, where a region of phase space is available where the potential is flat. For increasing $N_{\text{ext}}$ values, the overestimation of $P_A(\lambda_B|\lambda_A)$ decreases, where the REPPEXTIS simulations with $N_{\text{ext}} \geq 4$ are seen to align with the RETIS estimate. The relative errors for $P_A(\lambda_B|\lambda_A)$ for the RETIS and PPTIS simulations are 7.3 % and 8.7 %, respectively. There is not yet an implementation to extract error estimates for the REPPEXTIS simulations, which will be considered in future work to provide a thorough analysis.

For an $N_{\text{ext}}$ value of 4, the [5±4] ensemble samples paths that go directly from $\lambda_2$ to $\lambda_8$, where the cosine barrier $M$ must be felt. This does not mean, however, that the REPPEXTIS methodology with $N_{\text{ext}} \geq 4$ captures all memory within the $\lambda_2$ to $\lambda_8$ region for *all* systems. If metastable states (e.g. Gaussian wells) would be introduced at positions $M1$ and $M2$ (see Fig. 6.3), the [5±$N_{\text{ext}}$] paths would 'circle' around these states. This would result in chains of consecutive $\left[5^{\pm}\right]$RMR-$\left[6^{\pm}\right]$LML path segments at $M1$ and $M2$. After $N_{\text{ext}}/2$ such 'circles', the paths would then also tunnel from $M1$ to $M2$.

As such, $[i±N_{\text{ext}}]$ ensembles do not guarantee full memory retention over the entire $[\lambda_{i-N_{\text{ext}}-1}, \lambda_{i+N_{\text{ext}}+1}]$ interval. The only way to ensure such *spatial memory* over an interval $[\lambda_C,\lambda_D]$ is to define RETIS(-like) path ensembles that require paths to start from $\lambda_C$, and track progression along $\lambda$ until $\lambda_C$ or $\lambda_D$ is hit. REPPEXTIS offers increased spatial memory to its direct neighboring ensembles, (extensions of a $\left[i^{\pm}\right]$ segment must result in a $\left[(i-1)^{\pm}\right]$ or $\left[(i+1)^{\pm}\right]$ segment), but for larger $N_{\text{ext}}$ values, the additional memory is *temporal* in nature, as it is no longer known in which ensembles the path extensions will end up.

### 6.7.2 Rugged energy landscape

A two-dimensional rugged energy landscape (Fig. 6.5) was constructed by random deposition of 500 Gaussian wells. Harmonic walls ($k = 250$) similar to the shuffleboard potential were used to construct a square box ($[L_x, L_y] = [0.9, 0.9]$). The Gaussian wells were uniformly distributed in the domain $y \in [-0.225, 1.125], x \in [0.09, 0.81]$, with uniformly distributed heights in the domain $[-V_m, V_m] + V_m/4$ ($V_m = 0.075$), and uniformly distributed $x$- and $y$-widths in the domain $[L/3, L]$ ($L = 0.05$). States $A$ and $B$ were defined equivalently to the shuffleboard potential (cosine dips), with $x_A = 0.025$, $x_B = 0.875$, and $V_A = V_B = -0.15$. The overall slope $V_{\text{slope}}(x) = 0.2x$ was also
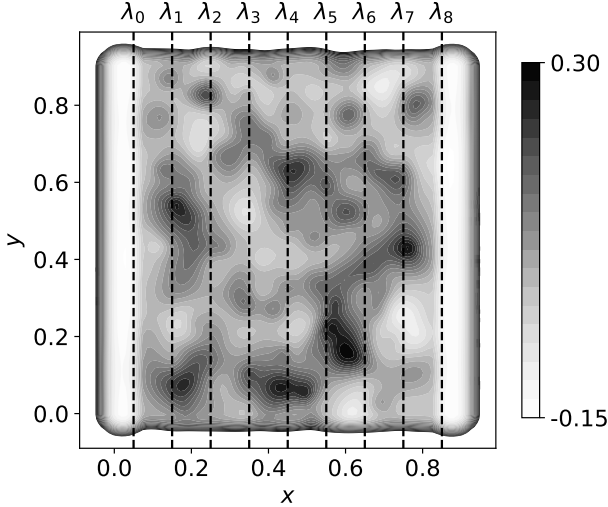
**Figure 6.5:** The rugged energy landscape. The Gaussian wells are shown as shades, where the color denotes the height of the well. The harmonic walls are shown as black lines.

added. The Langevin dynamics were integrated using the same reduced units as in the shuffleboard potential. Only the temperature was lowered, where the parameters used are: temperature $T = 0.04$, particle mass $m = 0.1$, friction coefficient $\gamma = 25$, and integration timestep $\Delta t = 0.002$.

Eight interfaces were positioned equidistantly, with $\lambda_A = 0.05$ and $\lambda_B = 0.85$. A RETIS simulation was performed as benchmark, where 75 % of the MC moves were shooting moves, and the other 25 % swapping moves. REPPEXTIS simulations were performed with $N_{\text{ext}} = 0$, 1, 2, 3 segment extensions, where strategy 3 (binomial extensions) was used for the shoot-and-extend moves. REPPEXTIS simulations with $N_{\text{ext}} = 1, 2, 3$ were also performed using $N_l = 4$ 'levels' of replicas, where each REPPEXTIS ensemble was thus included 4 times. As the permenant calculations proved to be computationally expensive for the $N_l = 4$ simulations ($4 \times 12 = 48$ ensembles), swapping moves did not consider all the path ensembles at once. Instead, two *swapping frames* were defined, where the set of 48 ensembles is partitioned into 3 or 4 sets, for which individual infinite swapping moves are performed. The first swapping frame partitions the ensembles according

to their level (4 sets of the usual 12 REPPEXTIS ensembles), while the second swapping frame partitioned the ensembles according to the $[i\pm N_{\mathrm{ext}}]$ ensemble indices. The first partition of frame 2 consists of the $[[0^-\pm N_{\mathrm{ext}}], \ldots, [2\pm N_{\mathrm{ext}}]]$ ensembles over all 4 levels, the second partition entailed the $[[3\pm N_{\mathrm{ext}}],\ldots,[6\pm N_{\mathrm{ext}}]]$ ensembles over all 4 levels, and the third partition entailed the $[[7\pm N_{\mathrm{ext}}],\ldots,[(N-1)^-\pm N_{\mathrm{ext}}]]$ ensembles over all 4 levels. At the start of an infinite swapping move, either of the two swapping frames is selected with equal probability.

The simulations ran for $\sim 10$ hours on a dell XPS 15 7590 laptop. Simulation cycle numbers (in thousands) were (45, 38, 31, 26) for the $N_{\mathrm{ext}}$ values (0, 1, 2, 3), respectively using one replica level ($N_l = 1$). For the $N_l = 4$ simulations, the cycle numbers were (in thousands) (8, 7, 6) for the $N_{\mathrm{ext}}$ values (1, 2, 3). As one cycle considers shooting moves in all of the path ensembles, a fair comparison of cycle counts should consider multiplication of the $N_l = 4$ cycle counts by 4. The RETIS simulation was performed for $10\,000$ cycles, resulting in a $P_A(\lambda_B|\lambda_A)$ estimate with a relative error of $37\,\%$. The PPTIS simulation ($N_{\mathrm{ext}} = 0$, $45\,000$ cycles) resulted in a relative error of $28\,\%$. The results are shown in Fig. 6.6. Apart from the PPTIS simulation, all REPPEXTIS simulations show good agreement with the RETIS estimate, indicating that memory effects are not significantly present. There is no significant difference between the multi-level and single-level replica simulations. Currently, an analysis framework to investigate the convergence rate is still lacking for the REPPEXTIS simulations, so a claim regarding improved convergence rate for multi-level replica simulations cannot be made. As increased replica exchange has consistently proven to be beneficial (TIS $\to$ RETIS $\to$ $\infty$RETIS for increased convergence rate), it is expected that future work will reveal a similar trend for REPPEXTIS.

## 6.8  Multidimensional approach

The application of REPPTIS to biological systems has raised the need for a multi-dimensional path ensemble definition. As the dimensionality increases, the amount of path ensembles grows (roughly) exponentially, where efficient sampling strategies are required to make the methodology feasible. The path extension strategy of REPPEXTIS could potentially be a powerful tool to allow efficient (infinite) replica exchange in such multi-dimensional spaces. In Fig. 6.7, a proposal for a two-dimensional implementation of the REPPEXTIS methodology is shown.
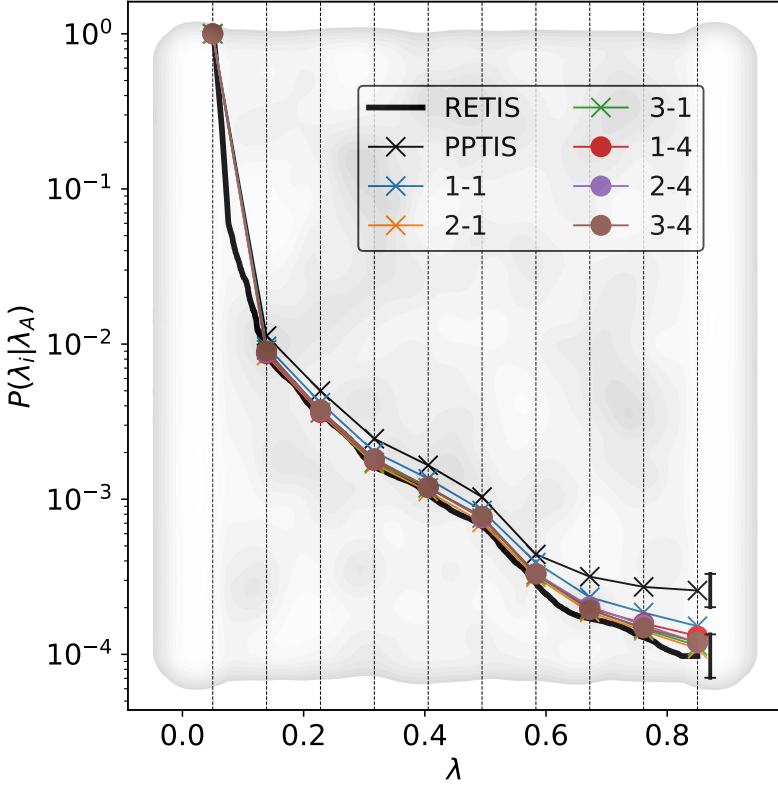
**Figure 6.6:** The crossing probability profiles $P_A(\lambda_B|\lambda_A)$ for the RETIS, PPTIS, and REPPEXTIS simulations. The rugged energy landscape is shown as a shade. The error bars correspond to the standard error of the RETIS and PPTIS $P_A(\lambda_B|\lambda_A)$ estimates. Labels of REPPEXTIS simulations denote the $N_{ext}$-$N_l$ values. X-marks represent the REPPEXTIS simulations using 1 level of replicas (conventional), while circles represent the REPPEXTIS simulations using 4 levels of replicas. An error estimate calculation for the REPPEXTIS simulations is not yet implemented.

## 6.9 Conclusion

In this chapter, a novel path sampling methodology, REPPEXTIS, was introduced. By direct sampling of chains of PPTIS path segments, the methodology showed promising results for memory enhancement (**Research Objective 1b**) in a shuffleboard potential. While it is expected that the infinite swapping formalism will result in increased ergodic sampling and improved convergence rates (**Research Objective 1a**), this has not yet been quantified with error estimates. Also the potential benefits of multiple replicas to increase sampling efficiency has not yet been investigated.
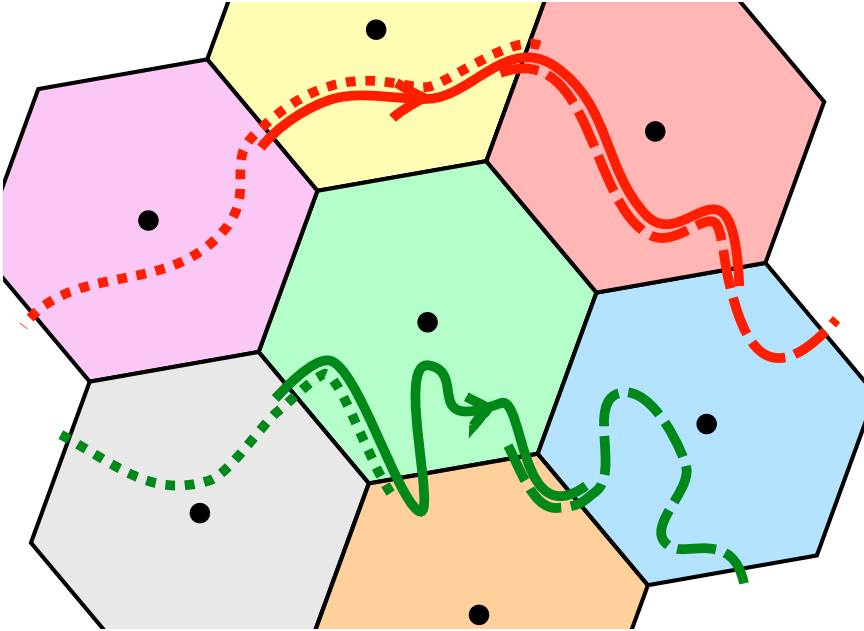
**Figure 6.7:**   A proposal for a two-dimensional implementation of the REPPEXTIS methodology. Interfaces are constructed using Voronoi tesselation on a selection of relevant CVs. In this example, hexagonally positioned interfaces are constructed. Each hexagon represents a multi-dimensional REPPEXTIS ensemble. A corresponding multi-dimensional PPTIS path segment starts at one of the hexagons edges $\lambda_1$, crosses *another* edge $\lambda_2$ *of the same hexagon* without recrossing $\lambda_1$, and then ends at *any* edge other than $\lambda_2$. As such, (L)eft, (M)iddle, and (R)ight parts of a trajectory can be identified just like done for the one-dimensional PPTIS MSM, where the (M)iddle part is the segment between the first and last crossing of $\lambda_2$. Path segments are then extended $N_{ext}$ times in randomly chosen time directions to create multi-dimensional PPTIS path chains. Two REPPEXTIS paths are shown, one of the green ensemble and one of the red ensemble. The bold paths are the original PPTIS paths in these ensembles, while the short-dashed and long-dashed paths represent extensions in the backward and forward time directions, respectively. The green REPPEXTIS path can be swapped into the gray, brown and blue ensembles. The red REPPEXTIS path can be swapped into the yellow, blue, and purple ensembles. The shooting move is to be adapted with a high acceptance formalism to select only phase points that are a part of its hexagon. This is effectively a Metropolis-based path sampling strategy of milestoning-like paths, where path memory can be extended to a user-specified degree ($N_{ext}$).

# 7

# MYELIN OXYGEN KINETICS

In this chapter, the kinetics of oxygen transport within and through myelin sheaths are studied. While the previous chapters explored path sampling methodologies to access longer timescales, a different approach is taken here that better suits the biological system at hand. Myelin sheaths can consist of up to more than a hundred phospholipid bilayers, where brute-force MD simulations of such large systems would be computationally demanding. Here, the periodicity of the system is exploited, where a diffusive network model for oxygen in one bilayer is constructed and subsequently expanded to represent an entire myelin sheath. In this model, the number of bilayers can be adjusted freely, allowing the study of **Research Question 2** pertaining to the effect of myelination thickness on oxygen storage and transport phenomena.

This chapter summarizes **paper IV** (steady-state oxygen phenomena) and **paper V** (time-dependent oxygen phenomena, **manuscript in preparation**). First, an introduction to the work is given, highlighting how the vital role of oxygen in maintaining neuronal function motivates the study of oxygen transport at the subcellular scale, where knowledge is lacking. Next, the methods used in the papers are shortly introduced, after which the results are summarized. The chapter ends with a discussion and conclusion.

## 7.1 Introduction

The brain, while only making up $2\,\%$ of total body weight, consumes over $20\,\%$ of the body's oxygen. [174], where neurons use the majority of the brain's energy to sustain their function, primarily through mitochondrial oxidative phosphorylation (OxPhos). Oxygen is vital as the final electron acceptor in the electron transport chain (ETC), imperative for maintaining the proton gradient necessary for ATP synthesis [174, 175, 212]. Cytochrome $c$ oxidase (COX) is the most crucial enzyme for molecular oxygen in the brain, and its high affinity for oxygen ensures undiminished activity even at low oxygen levels [213–215]. The pathways of oxygen towards COX are complex, consisting of many energetic and diffusive barriers, where the impact of subcellular structures on oxygen transport phenomena remains largely unknown. In this work, the subcellular structure of myelin sheaths is considered.

## 7.2 Oxygen storage in myelin

Myelin sheaths, composed of up to 100 tightly packed phospholipid bilayers, play a crucial role in insulating axons and facilitating the rapid transmission of electrical signals through saltatory conduction [176, 216, 217]. Despite the presence of tightly packed polar phospholipid headgroups, which pose both energetic and diffusive barriers, the overall membrane oxygen permeability remains relatively high. This phenomenon aligns with Overton's rule (high permeability due to high lipid solubility), and MD studies have further supported this by demonstrating efficient oxygen diffusion through the lipid bilayer [218–221]. Total oxygen solubility is 3 to 5 times larger in phospholipid bilayers than in water, reaching upwards of 10 times in the hydrophobic mid-plane of the bilayer [222–226].

In **paper IV**, oxygen storage in stacks of phospholipid bilayers is investigated. To this end, a diffusive network model was constructed, based on previously published MD simulations of oxygen molecules in a 1-palmitoyl-2-oleoyl-sn-glycero-3-phosphocholine (POPC) bilayer [219]. The free energy profile $F$ and the diffusivity profile $D$ were extracted from these simulations using the Bayesian inference approach developed by Hummer [227]. Both a snapshot of the simulation box and the extracted $F$ and $D$ profiles are shown in Fig. 7.1. From these profiles, a corresponding rate matrix can be constructed. To create a diffusive rate model for myelin, the $F$ and $D$ profiles of one bilayer were cut and subsequently appended $M$ times to construct
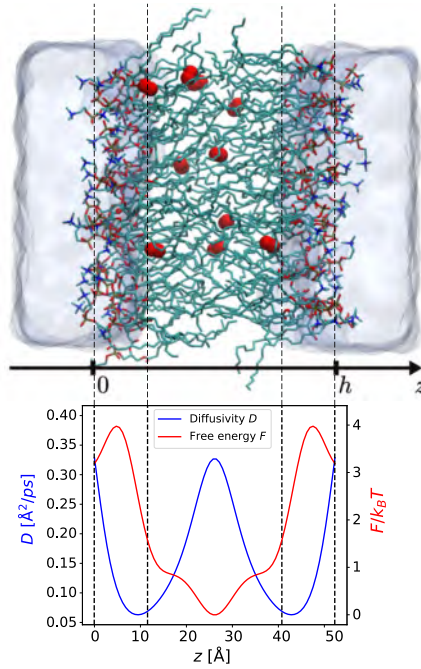
**Figure 7.1:** A snapshot of the simulation box (top) and the extracted $F$ and $D$ profiles (bottom). The coordinate $z$ is perpendicular to the bilayer plane. The bilayer is contained in the interval $z \in [0, h]$, where the bilayer thickness $h$ is approximately $52\,\text{Å}$.

profiles corresponding to a stack of $M$ phospholipid bilayers. The rate matrix was then manipulated to remove the periodicity of the system, after which the introduction of constant concentration boundary conditions allowed the study of steady-state oxygen transport.

It was found that the POPC membrane stores about 8 times more oxygen than pure water. Stacking bilayers, as in myelin sheaths, enhances oxygen storage capacity but with diminishing returns after a certain number of layers. The presence of additional water between stacked bilayers, as sometimes observed in cancer cells, negatively impacts the oxygen storage enhancement. Due to the large presence of myelin in the brain, and due to their close proximity to the oxygen consuming axonal mitochondria, the hypothesis of myelin acting as an oxygen reservoir arises. To further investigate this hypothesis, the time-dependent kinetics of oxygen transport through myelin sheaths were examined.

## 7.3 Oxygen kinetics in myelin

A large gap between the two smallest eigenvalues of a single bilayer rate matrix suggested that oxygen kinetics followed first order kinetics. In **paper V**, time-dependent boundary conditions were introduced to the rate models, where the (un)loading of oxygen to a bilayer was indeed found to be accurately described by a single exponential decay (or growth) function. The (R)esistance of a bilayer is a well known property, and is given by the inverse of the bilayer permeability. A (C)apacity of a bilayer was defined as the total oxygen storage at equilibrium conditions, thus denoting its capacity for oxygen storage. Using this, a simple and intuitive RC circuit analogy was constructed that accurately captures the time-dependent oxygen storage of a bilayer. In the supplementary information of **paper V**, it is described how the general shape of the $F$ and $D$ profiles lend a derivation of such an RC circuit directly from the Smoluchowski equation, and that the model is thus assumed to be applicable for other bilayers and/or (non-polar) permeants as well.

Using this RC circuit analogy, a myelin sheath was then represented by a ladder circuit of $N$ RC circuit elements (Fig. 7.2). It was
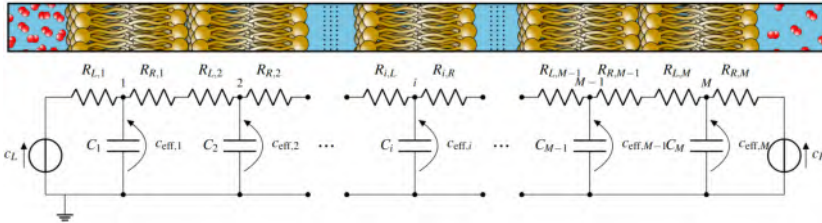


**Figure 7.2:** A myelin sheath represented by a ladder circuit of $N$ RC circuit elements. Notations are explained in-depth in **paper V**.

found that the characteristic time-constant $\tau$ of oxygen storage increased quadratically with the number of bilayers. Concretely, this ranges from ∼30 ns for one POPC bilayer up to ∼506 µs for 200 POPC bilayers, illustrating the very large increase in oxygen transport time scales when stacking bilayers. Is half a millisecond relevant? To answer this, some numbers are due.

As oxygen consumption is central to energy production, it is often used to probe neuronal activity. The blood oxygen level dependent (BOLD) signal detected in functional magnetic resonance imaging (fMRI) often displays an 'initial dip', which has been linked to a quick increase in the cerebral metabolic rate of oxygen ($CMRO_2$) prior to

the cerebral blood flow (CBF) increase [228, 229]. The vascular response (CBF increase) occurs $\sim 500\,\mathrm{ms}$ after the neural stimulus, and quickly thereafter dominates the increase of $CMRO_2$ [230, 231]. Directly probing $P_{O_2}$ in capillaries using two-photon lifetime microscopy (2PLM) also displays the initial dip in tissue $P_{O_2}$ following neuronal activation [232, 233]. This dip is observed within $\sim 100\,\mathrm{ms}$ following synaptic activation, and preceding functional hyperemia by $\sim 1\,\mathrm{s}$ [232]. The shape and duration of the dip are, however, unpredictable, as they depend highly on local parameters [232, 234]. It was recently shown how this dip in $P_{O_2}$ is causal to an increased local blood flow response (capillary hyperemia), even in the absence of neuronal stimuli [235].

While most of the listed timescales are two orders of magnitudes larger than one millisecond, these experiments either measure (a) averaged effects and (b) measure $P_{O_2}$ in or near capillaries, and not in the myelin that is in much closer proximity to the axonal mitochondria (as it enwraps them). If a significant decrease in $P_{O_2}$ is observed near the capillary as soon as $100\,\mathrm{ms}$ after neuronal activation, it is plausible to assume that the $CMRO_2$ increase occurs (much) faster. Direct measurements of $CMRO_2$ at the singular axon level were, however, not (yet) found.

To further bridge the gap between theoretical and experimental timescales, the RC ladder circuit was extended to include (larger) circuit elements that modelled the axonal cytosolic solvent and extracellular matrix, modeling the oxygen transport from capillary to mitochondria as visualized in Fig. 7.3. As noted before, neuronal activity is met with an increase in oxygen consumption. A $\sim$ 5-fold increase in the oxygen consumption rate (OCR) has been reported in the hippocampal CA3 network during gamma oscillations as compared to absence of spiking (equivalently, a $\sim$ 2-fold increase as compared to spontaneous activity) [236]. The RC ladder model was adapted to include time-dependent flux boundary conditions, where the efflux at the axonal mitochondrial represented the OCR. Neuronal activity was then modelled by increasing the OCR $f$-fold over a characteristic time $\tau_{\mathrm{OCR}}$. The impact of myelination was then studied by calculating the time for which the increased OCR could be sustained (i.e. until the axonal oxygen concentration becomes zero). Various settings were tested, where the results showed that sustainment times increased for increased myelination. The axonal oxygen concentration at rest conditions and the value of $\tau_{\mathrm{OCR}}$ were found to have a more significant impact on the sustainment time. Interestingly, most of the scenarios
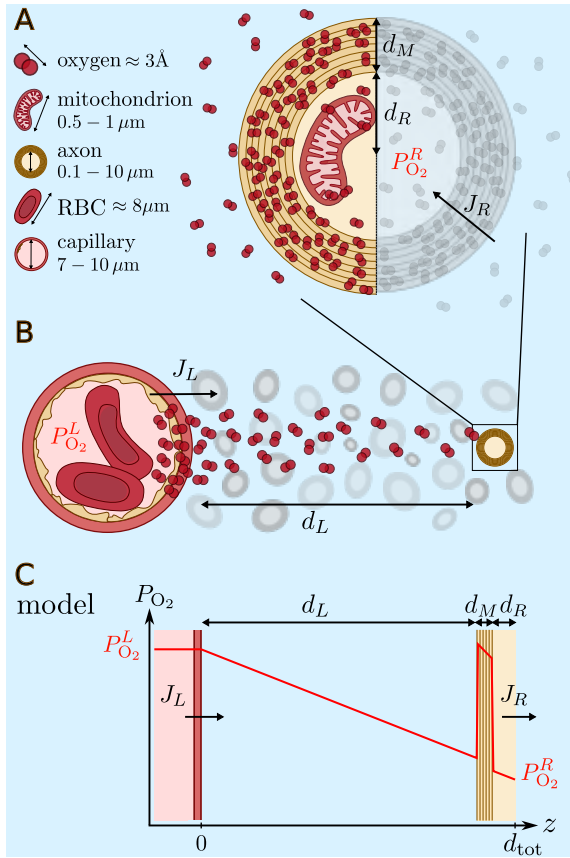
**Figure 7.3:** (A). Sketch of myelinated axon. (B). Sketch of capillary and myelinated axons. (C). Model of transport of oxygen from the capillary to the myelinated axon. Indicated distances: $d_R$ is the radius of the axon, $d_M$ is the thickness of $M$ bilayers in the myelin sheath, $d_L$ is the distance from the myelin sheath to the capillary. Indicated fluxes: $J_L$ (left) is the $O_2$ flux flowing from capillary, $J_R$ (right) is the $O_2$ flux reaching axonal mitochondria.

considered could *not* sustain the increased OCR for much longer than 100 ms, underscoring the functional role of the hyperemic blood flow response that quickly follows thereafter.

The RC ladder model to study the neuronal set-up of Fig. 7.3 is a large simplification of the biological system, where myelin sheaths are represented as infinitely extended sheets of phospholipid bilayers. While a radial solution to the set-up has been explored, it would be most interesting to create a model containing a distribution of cylindrical myelinated axons and capillaries with the inclusion of more realistic source and sink boundary conditions.

## 7.4 Conclusion

The oxygen storage and transport phenomena in myelin sheaths were studied using a diffusive network model and an RC ladder circuit analogy. The results showed that the oxygen storage capacity of a myelin sheath increases with the number of bilayers, but with diminishing returns. The time-constant of oxygen storage in a myelin sheath was found to increase quadratically with the number of bilayers, illustrating the large increase in oxygen transport time scales when stacking bilayers. The RC ladder circuit analogy was extended to include larger circuit elements that modelled the axonal cytosolic solvent and extracellular matrix, modeling the oxygen transport from capillary to mitochondria. The impact of myelination on the sustainment of increased oxygen consumption rate was studied, where the results showed that sustainment times increased for increased myelination. The results underscore the functional role of the hyperemic blood flow response that quickly follows neuronal activation.

To conclude **Research Question 2**, the effect of myelination on time-(in)dependent oxygen storage was successfully characterized, while its impact on oxygen transport phenomena at larger scales remains to be further explored.

# CONCLUSIONS AND OUTLOOK

Biological processes occur across a wide range of timescales, many of which include (configurational) transitions that are beyond the reach of conventional MD simulations. The kinetic analysis of such rare and/or slow transitions exacerbates this timescale gap, as observing only a few transitions is insufficient for statistically significant estimations. This work focused on developing and applying computational methodologies to extract the kinetics of long-timescale biological processes at the molecular scale. The first six chapters addressed the first Research Question: *Can the PPTIS methodology be improved to provide reliable kinetic analysis of long-timescale biological processes?* To this end, the limitations of the PPTIS methodology were identified and addressed through four intermediate Research Objectives. These limitations included (1a) the lack of ensemble information exchange, (1b) the lack of path memory, (1c) the inability to extract time-dependent properties, and (1d) the absence of an implementation to apply PPTIS using MD engines.

The paths of PPTIS ensembles are confined to a small interval along a chosen order parameter $\lambda$. Energy barriers orthogonal to $\lambda$ can hinder these paths from sampling relevant regions of phase space. A strategy to alleviate this limitation was already proposed in the original PPTIS paper by means of a replica exchange (path swapping) formalism. In this approach, paths from neighboring ensembles are extended into each other's $\lambda$ interval and then exchanged. This

replica exchange formalism was implemented, where the methodology was coined replica exchange PPTIS (REPPTIS). The increase in ergodic sampling was demonstrated in **paper I**, where an improved permeability estimation was obtained for a low-dimensional maze potential. It was shown that orthogonal barriers could trap PPTIS paths within a specific region of phase space, while the replica exchange moves enabled the sampling of paths across these barriers Furthermore, the PPTIS and REPPTIS methodologies were integrated into the PyRETIS software package in **paper II**, facilitating their application with common MD engines such as Gromacs, LAMMPS, and OpenMM. Thus, objectives (1a) and (1d) were successfully achieved.

By viewing long MD trajectories as a chain of overlapping PPTIS paths, a Markov state model (MSM) was constructed in **paper III** (manuscript in preparation) to estimate the length of trajectories that are much longer than the short PPTIS paths. In this MSM, each of the four possible path types of every PPTIS ensemble is treated as a state, with transition probabilities defined by the local crossing probabilities of the ensembles. This allows for the estimation of the probability of longer trajectories and their corresponding path lengths through a random walk reconstruction in the MSM. Path lengths are directly related to the *path durations* by a simple timestep scaling factor. By calculating the average path durations of both $[0^+]$ and $[0^-]$ RETIS paths, the flux term entering rate equations can be estimated. As such, objective (1c) was achieved.

The REPPTIS methodology was applied to protein-drug complexes, where the residence time of a drug molecule to its target protein is a crucial estimator for *in vivo* drug efficacy. In **paper III**, the unbinding kinetics of the benzamidine molecule from the trypsin kinase protein was studied. The dissociation rate was underestimated by 4 orders of magnitude compared to experimental data and other computational studies. This discrepancy was presumably due to a poor initialization of the PPTIS paths, where the initial trajectory generated by steered MD was not representative of the dominant reactive pathway. Replica exchange moves were infrequently performed in the ensembles located near the bound state of the benzamidine molecule, suggesting that paths were unable to escape the initially biased region of phase space towards more relevant regions of path space. A long MD simulation revealed that this underestimation was already present in the crossing probabilities related to these ensembles closely located to the bound state. As such, the combination of this long simulation and the PPTIS methodology was still successful in estimating

the dissociation rate. However, it becomes evident that REPPTIS becomes dependent on the initial path generation, and that the replica exchange move is not a guarantee for ergodic sampling.

In **Chapter V**, the dissociation of imatinib from the ABL kinase protein and mutated ABL variants is investigated. Despite a more careful construction of the initial paths, the REPPTIS simulations did not converge to reliable rate estimates. Challenges in sampling arose due to metastable states that are separated by energy barriers orthogonal to the $\lambda$ parameter, creating hidden timescale separations that are not easily overcome by the REPPTIS methodology. An $\infty$RETIS simulation was then performed to investigate dissociation close to the bound state. Even in close proximity to the deep binding pocket, a large timescale separation persisted, indicating the presence of metastable states. This study concludes that the path sampling methodologies faced significant limitations in accurately capturing the dissociation process due to the complexity of the energy landscape not being adequately represented by a one-dimensional order parameter. The rugged energy landscape of many biological systems is expected to be similar to that of the ABL-imatinib complexes considered here. This highlights the need for the REPPTIS methodology to better handle the presence of metastable states, their orthogonal energy barriers, and their associated hidden timescale separations. A major advancement in this direction will involve developing path ensembles using multi-dimensional order parameters.

In **Chapter VI**, the novel path sampling methodology REPPEXTIS was developed to enhance path memory (objective 1b) and further increase ergodic sampling of path space (objective 1a). REPPEXTIS paths inherently represent continuous chains of PPTIS path segments due to path extensions performed subsequent to a successful shooting move. Consequently, states in the REPPEXTIS MSM are defined as all possible $N_{\text{ext}} + 1$ PPTIS segment chains of the PPTIS MSM, thereby increasing path memory to a user-specified degree. REPPEXTIS paths belong to multiple path ensembles, allowing swap moves to be performed without the need for MD integration. As such, the infinite swapping formalism that was recently introduced in $\infty$RETIS could be adapted and incorporated. Additionally, REPPEXTIS introduces the concept of multiple replicas of the same path ensemble to maximize sampling and information exchange in difficult-to-sample regions. The methodology was tested on two simple toy systems with promising results, but its full potential remains to be explored.

To conclude Research Question 1: the REPPTIS methodology has significantly improved upon the PPTIS methodology, but it still faces a strong limitation due to its reliance on a one-dimensional order parameter. In its current form, the REPPTIS methodology is not fully suited to study biological systems where the reaction process is not well-described by a one-dimensional order parameter. A clear next step involves the development of multi-dimensional path ensembles, allowing importance sampling to focus on path segments that cross barriers in multiple directions. A proposal for a multi-dimensional implementation of the REPPEXTIS methodology was given at the end of Chapter 6, where the infinite swapping formalism due to the path extensions may provide adequate sampling of the large number of path ensembles that would arise when multiple dimensions are considered.

To tackle the second Research Question, *How does myelination affect the storage of oxygen and its transport to axons?*, a diffusive network model was constructed to study the slow oxygen kinetics within myelin sheaths. In **paper IV**, a trick was used to artificially enlarge MD simulation data of oxygen in one phospholipid bilayer. A myelin sheath was represented as a stack of $N$ phospholipid bilayers, where the free energy and diffusivity profiles of one bilayer were simply concatenated. It was demonstrated that stacked phospholipid membranes increase oxygen storage capacity, suggesting that myelin may serve as an oxygen reservoir for the nearby oxygen-consuming cytochrome c oxidase in (axonal) mitochondria. It is further shown that this effect levels off for a large number of bilayers, and that additional water between layers, as seen in some cancer cells, diminishes oxygen storage enhancement.

In **paper V** (manuscript in preparation), the model was extended with time-dependent boundary conditions to characterize the loading and unloading effects of oxygen in myelin sheaths. It was demonstrated that oxygen (un)loading in one phospholipid bilayer follows first-order kinetics. By using the membrane resistance (R, inverse permeability) and defining a capacitance (C, capacity for oxygen storage) a simple RC circuit analogy was built that accurately and intuitively captures time-dependent oxygen storage and release in one bilayer. A ladder circuit of RC elements was then constructed to describe oxygen transport through myelin. Stacking multiple bilayers increases the total oxygen storage capacity linearly. However, the time constant for oxygen storage increases quadratically with the number of bilayers, indicating a slower response to changes in oxygen concentration as the

myelin sheath becomes thicker. By embedding the ladder circuit in extracellular and axonal solvent compartments, oxygen transport from capillary to axonal mitochondria was investigated. Neuronal activity, characterized by an increase in oxygen consumption rate (OCR), was modelled by increasing the oxygen efflux boundary condition at the axonal compartment. The effect of myelination was then investigated by calculating the duration for which the increased OCR could be sustained (i.e. the time until the axonal oxygen concentration becomes zero). It is shown that higher levels of myelination increase the time for which increased neuronal activity can be sustained. This magnitude of this effect is, however, highly dependent on (a) the onset time of the OCR increase, and (b) the axonal oxygen concentration at rest conditions. Multiple set-ups were considered, and interestingly none of the scenarios could sustain the OCR increase for much longer than one hundred microseconds. This result suggests that neuronal activity may be met with hypoxia in the axonal compartment, and that the dominant 'overshooting' vascular response is necessary to restore oxygen levels.

## 8.1 PERSPECTIVES

At the end of a long journey, a moment of reflection is due. In this final section, perspectives are given on the journey towards the development of the introduced methodologies, highlighting their achievements and limitations, and determining directions for future work.

Science is a journey of discovery, often involving more than just finding answers to predefined questions; it also entails identifying the right questions to pursue. This was true for the development of the REPPTIS methodology, where initially RETIS was selected as a tool to study the effect of mutations on drug binding kinetics, contributing to the field of personalized medicine. Due to the presence of long-lived metastable states, the application of RETIS was not feasible, seeding the development of the REPPTIS methodology. While a correct method-development protocol was followed - testing the methodology first on low-dimensional toy systems - the limitations of REPPTIS might have been more easily discovered by further testing on less complex biological systems than the ABL-imatinib complexes. The application of REPPTIS to the less complex trypsin-benzamidine system was initiated considerably later than the ABL-imatinib systems. In retrospect, it would have been more fruitful if the systems were treated in reverse order, as the limitations were already present, albeit to a lesser extent, in the trypsin-benzamidine system.

While limitations of the REPPTIS methodology were discovered, it still serves as a powerful tool to study biological systems for which the reaction process is both rare (involving high energy barriers) and slow (involving metastable states). However, its application is limited to reactions whose transition pathways are well-described by a one-dimensional order parameter. The construction of initial paths for the REPPTIS ensembles is therefore crucial, as it is likely that user-defined order parameter will deviate from ideal reaction coordinates, such as the committor function. Especially if orthogonal energy barriers separate metastable states, an initial path that deviates from the dominant reactive pathway can result in paths being stuck in 'orthogonally locked' unfavorable regions of phase and path space. Careful construction of such initial paths is therefore required, where the slowly performed pulling simulations used in this work are thought to be insufficient for proper path initializations.

Many of the aforementioned issues of multiple reactive pathways, orthogonal energy barriers and their hidden timescale separations could be addressed by effectively eliminating their orthogonality. A most promising direction for future work therefore lies in the development of a multi-dimensional extension of the memory-truncated TIS-based methodologies like PPTIS, REPPTIS, and REPPEXTIS. A proposal for a multi-dimensional implementation of the REPPEXTIS methodology was given at the end of Chapter 6, where the combination of extensions and infinite swapping may enable sufficient sampling of the larger number of path ensembles involved. A careful analysis of the energy landscape could then be used to extract the dominant collective variables of the reaction process, which are to be used as $\lambda$ parameters for the multi-dimensional REPPEXTIS simulation.

AI has assisted the non-specific search for dominant collective variables, where for example VAMPnets [157] saw success in the Markov state modelling approach to kinetics (see Sec.1.3). However, the predictive capacity of AI methodologies applied to kinetics is not yet at the level of thermodynamical inquiries. During the four year-long duration of this PhD, the scientific community saw the release of DeepMind's AlphaFold [237, 238], and OpenAI's GPT large language models [239]. An essential element to AlphaFold's success in protein structure prediction (and other neural nets that predict protein-ligand interfacing) was the availability of extensive data on proteins, small drug molecules, and their complexes (ChEMBL [240], PubChem [241], BindingDB [242], DrugBank [243], RCSB PDB [244], UniProt [245],

etc.). While efforts such as the 'kinetics for drug design' consortium (K4DD [246]) have been made to establish a database targeting protein-drug (un)binding kinetics, it will likely take time before database-sizes are sufficiently large for neural nets to provide on-the-spot predictions on rare-event time scales (such as '*what is the dissociation time of drug D from protein P?*'). AI is, however, paving its way in the kinetics field, by aiding with collective variables (e.g. VAMPnets), or biasing TPS towards to the generation of more reactive pathways (e.g. machine-guided TPS [247]). It will be particularly interesting to see how AI-guided path sampling is adopted from TPS to TIS-based methods, especially when combined with machine-guided selection of important order parameters.

With respect to second part of this dissertation (i.e. the study of oxygen storage and transport in myelin), the main effort of future work will lie in (1) collaboration to obtain experimental data, and (2) a structurally more complex model of myelin. Earlier discussions with experimental colleagues have, however, indicated that probing oxygen kinetics at the scale of our models may not be feasible. Therefore, a larger and more detailed model should be constructed from which observables at larger scales can be derived. Such a model should would ideally involve a distribution of cylindrical axons and capillaries, with the inclusion of more realistic source and sink boundary conditions to model the increased oxygen demand following neuronal activity.

# II

---

# Part II – Papers

---

# 9

# PAPER I (PUBLISHED): PATH SAMPLING WITH MEMORY REDUCTION AND REPLICA EXCHANGE TO REACH LONG PERMEATION TIMESCALES

W. Vervust constructed the model systems, performed the (RE)PPTIS simulations, and contributed to data analysis and writing of the manuscript.

Biophysical Journal
## Article

Biophysical Society

# Path sampling with memory reduction and replica exchange to reach long permeation timescales

Wouter Vervust,[1] Daniel T. Zhang,[2] Titus S. van Erp,[2] and An Ghysels[1,*]

[1]IBiTech – Biommeda Research Group, Faculty of Engineering and Architecture, Ghent University, Gent, Belgium and [2]Department of Chemistry, Norwegian University of Science and Technology, Trondheim, Norway

ABSTRACT   Assessing kinetics in biological processes with molecular dynamics simulations remains a computational and conceptual challenge, given the large time and length scales involved. For kinetic transport of biochemical compounds or drug molecules, the permeability through the phospholipid membranes is a key kinetic property, but long timescales are hindering the accurate computation. Technological advances in high-performance computing therefore need to be accompanied by theoretical and methodological developments. In this contribution, the replica exchange transition interface sampling (RETIS) methodology is shown to give perspective toward observing longer permeation pathways. It is first reviewed how RETIS, a path-sampling methodology that gives in principle exact kinetics, can be used to compute membrane permeability. Next, recent and current developments in three RETIS aspects are discussed: several new Monte Carlo moves in the path-sampling algorithm, memory reduction by reducing pathlengths, and exploitation of parallel computing with CPU-imbalanced replicas. Finally, the memory reduction presenting a new replica exchange implementation, coined REPPTIS, is showcased with a permeant needing to pass a membrane with two permeation channels, either representing an entropic or energetic barrier. The REPPTIS results showed clearly that inclusion of some memory and enhancing ergodic sampling via replica exchange moves are both necessary to obtain correct permeability estimates. In an additional example, ibuprofen permeation through a dipalmitoylphosphatidylcholine membrane was modeled. REPPTIS succeeded in estimating the permeability of this amphiphilic drug molecule with metastable states along the permeation pathway. In conclusion, the presented methodological advances allow for deeper insight into membrane biophysics even if the pathways are slow, as RETIS and REPPTIS push the permeability calculations to longer timescales.

SIGNIFICANCE   Permeability is a key kinetic property for membranes. Simulating permeation events at the molecular scale is very valuable for kinetic modeling, but permeation timescales are often prohibitively long to be simulated with present-day computational resources. Here, we show how the RETIS path-sampling method can give the exact kinetics of permeation, and how its efficiency is aided by recent developments. Moreover, we present a newly implemented REPPTIS method that approximates the kinetics by truncating memory. The REPPTIS method is promising for permeation simulations with high efficiency and accuracy that might not be easily achieved by any other method.

## INTRODUCTION

Biological membranes are responsible for compartmentalization in cells and organelles. Their permeability is a key characteristic of the transport kinetics of chemicals and nutrients, peptide-membrane interactions, or drug delivery of nanocarriers (1–7). Molecular dynamics (MD) simulations are a computational tool that aid the understanding of the

biophysical mechanisms playing at the molecular scale. Unfortunately, membrane permeability simulations are often computationally very demanding. The study of permeation events requires long timescales when the permeation is a slow and/or rare event, in addition to the need for fairly large simulation boxes usually comprised of thousands of particles. Permeability methods, such as the counting method (8–10) and the inhomogeneous solubility-diffusion model (11–13), are hence hindered by poor statistics. The latter approach can be combined with methods such as umbrella sampling (14) or adaptive biasing force (15,16) to obtain the free energy profile more efficiently. However, the possible presence of hysteresis and parallel reaction

channels can still sabotage an accurate description of the dynamics.

In recent work by some of the authors, an algorithm was proposed to evaluate the permeability with the path-sampling methodology, which realizes a speed-up of several orders of magnitude when the permeation event is rare (17). Path sampling, and in particular transition interface sampling (TIS) (18), achieves this speed-up by Monte Carlo (MC) sampling of path ensembles that are populated mostly with reactive paths or paths that make substantial progression along the reaction coordinate before returning to the reactant state. It was derived how the permeability can be obtained from the replica exchange transition interface sampling (RETIS) method (17). The method does not need a diffusive assumption, and the kinetics are thus exact.

Other notable methods that try to determine dynamic quantities, faster than MD, are milestoning (19) and forward flux sampling (FFS) (20). The latter is based on the same theoretical foundations as TIS, but instead of Metropolis MC sampling (21), it is based on splitting. In this class of methods, phase points of trajectories far up the barrier are used to launch multiple trajectories that deviate from the original due to the stochastic nature of the dynamics. Some of these trajectories reach further and deliver new phase points for launching the next set of trajectories and so on. Although the FFS method has the advantage over TIS and RETIS that it can be applied for nonequilibrium dynamics, the splitting technique has as a major disadvantage in that the final transition trajectories can be highly correlated and that possibly important rare initial conditions with a high reaction probability are missed (22). Milestoning, on the other hand, has the advantage that the statistics of long transition paths is obtained via the transition probabilities of much shorter paths connecting so-called milestones. The method bears resemblance to the simultaneously introduced TIS variation called partial path TIS (PPTIS) (23). Unlike the TIS, RETIS, and FFS methodologies, both PPTIS and milestoning are generally not exact as they rely on a memory loss assumption (Markovian approximation). Only in the hypothetical case that the milestones/interfaces are identical to isocommittor surfaces do these two methods become exact (24). In any other case, the inclusion of some memory can improve the accuracy of the method, which we discuss later.

Of the above methods, milestoning has been most commonly used to calculate permeability through membranes via slightly different theoretical approaches based on Markovian approximations (25–28). We recently derived how the permeability can be computed from a RETIS simulation that is not based on a disturbance out of equilibrium or Markovian assumptions (17), but, instead, it reproduces exact results identical to the counting method in a hypothetically long equilibrium MD run.

RETIS has been shown to provide the permeability for a range of toy systems (17). Moreover, RETIS was successfully

used in a realistic simulation setup to compute the permeation rate of oxygen molecules through a 1-palmitoyl-2-oleoyl-*sn*-glycero-3-phosphocholine membrane (29). Nevertheless, for systems with slow permeation events, where the permeation trajectories sampled in the path ensembles occasionally are very long, the RETIS simulation can still be computationally demanding. To reach such long timescales, further methodological improvements are needed to improve sampling and efficiency.

This paper starts with a short review on how the permeability formula is derived from RETIS. Next, we give an outlook on how the simulations at these long timescales can be made more efficient. For permeation events with very long pathlengths, a lower computational cost may be obtained with a reduction of memory similar to milestoning and PPTIS. Specifically, we extend the PPTIS method with a replica exchange move between path ensembles to enhance nonlocal sampling. Other improvements for the sampling efficiency are reviewed as well, such as the use of special path-generating MC moves (17,30,31) and a new way to run replica exchange simulations with cost-unbalanced replicas with an infinite swap frequency (32).

The set of these methodological advances in three areas, i.e., memory reduction, new Monte Carlo moves, and parallel computing, will push permeability calculations to longer timescales. Specifically, the methodology that shortens memory is illustrated in the results using two example systems. In the first system, a permeant needs to pass through a maze-like membrane, choosing between a pathway with an entropic barrier and another pathway with an energetic barrier. The role of memory for PPTIS is challenged in this setup, as the memory loss might induce an overestimation or underestimation of the permeability compared with RETIS. In the second example, ibuprofen permeation through a bilayer is modeled. This drug molecule is amphiphilic, with a polar part (carboxylic group) and apolar part. The polar group creates metastable states in ibuprofen's permeation pathway, which can cause the trajectories to be trapped, making it an excellent application of memory reduction to keep the trajectory lengths short. The last section summarizes the conclusions.

## METHODS

### Permeability from RETIS

A RETIS simulation requires the definition of an order parameter $\lambda$ and $n + 1$ interfaces $\lambda_A = \lambda_0 < \ldots < \lambda_i < \lambda_n = \lambda_B$ to describe the progression of the reaction. For a permeability calculation, the order parameter is simply the coordinate of a specific "target" permeant, orthogonal to the membrane plane, so $\lambda = z_t$ in Fig. 1. In curved membranes, such as liposomes, $\lambda$ could be chosen as the radial distance (33). In this subsection, a short overview of the permeability derivation is given to introduce the quantities that are needed to compute
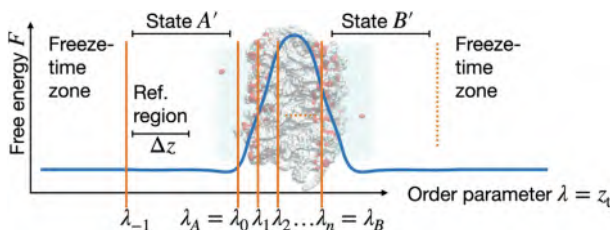
Vervust et al.

FIGURE 1 The membrane region between state A and state B forms a free energy barrier. Order parameter $\lambda$ is the $z_t$ coordinate of the target permeant normal to the membrane. RETIS interfaces $\lambda_A = \lambda_0, ..., \lambda_n = \lambda_B$ are indicated with orange lines. Additional interface $\lambda_{-1}$ reduces region A to region A′ by freezing the time. Reference interval $\Delta z$ is used to measure $\tau_{\text{ref},[0^-]}$ in Eq. 8 or $(\rho_{\text{ref}})_{A'}$ in Eq. 9. To see this figure in color, go online.

the permeability for the numerical applications in the next subsection.

A RETIS simulation consists of a series of path-sampling simulations (34), each employing a different path ensemble. There are $n + 1$ path ensembles called $[0^-]$, $[0^+]$, $[1^+]$, $[2^+]$, ..., $[(n-1)^+]$. The $[0^-]$ ensemble explores the reactant well and consists of paths starting and ending at $\lambda_0$ with all other path frames being at the left side of this interface ($< \lambda_0$). All other path ensembles $[i^+]$ with $0 \leq i < n$ explore the barrier region and consist of paths starting at $\lambda_0$, ending at either $\lambda_0$ or $\lambda_n$, and crossing interface $\lambda_i$ at least once. Since the path-sampling algorithm is based on MC moves obeying detailed balance, the set of trajectories that is being sampled in each path ensemble is statistically equivalent to a set of trajectories one would get by cutting out the relevant segments from an infinitely long MD trajectory.

Based on the results of the path ensembles, the rate can be computed as

$$k = f_A P_A(\lambda_B | \lambda_A) \qquad (1)$$

where $f_A$ is the conditional flux through $\lambda_A = \lambda_0$ and $P_A(\lambda_B | \lambda_A)$ is the overall crossing probability, the chance that $\lambda_B$ will be crossed after a positive crossing with $\lambda_A$ without recrossing $\lambda_A$. The rate in Eq. 1 can be related to a frequency of transitions between history-dependent states, called *overall states*, $\mathscr{A}$ and $\mathscr{B}$. These overall states differ from the stable states $A$ and $B$, which are defined as the phase space regions being at the left of $\lambda_A$ and at the right of $\lambda_B$, respectively. The overall state $\mathscr{A}$ also includes ("trajectory-") phase points that are in between $\lambda_A$ and $\lambda_B$, which were more recent in $A$ than in $B$. Likewise, overall state $\mathscr{B}$ includes points in between $\lambda_A$ and $\lambda_B$, which originate from paths that were more recently in $B$ than in $A$. Thus, the overall states ($\mathscr{A}, \mathscr{B}$) are larger than the stable states ($A, B$). The rate constant $k$ is then defined as the number of $\mathscr{A} \to \mathscr{B}$ transitions in a hypothetical infinitely long equilibrium MD run divided by the time spent in $\mathscr{A}$.

The flux $f_A$ in Eq. 1 is the frequency of positive crossing with $\lambda_A$ under the condition that the system is in overall $\mathscr{A}$. In rare events, however, the more complicated term to compute in Eq. 1 is the overall crossing probability $P_A(\lambda_B | \lambda_A)$ since it usually is an extremely small number.

In TIS, RETIS, and also FFS, it is computed from the following product expression:

$$P_A(\lambda_B | \lambda_A) = P_A(\lambda_n | \lambda_0) = \prod_{i=0}^{n-1} P_A(\lambda_{i+1} | \lambda_i) \qquad (2)$$

where $P_A(\lambda_{i+1} | \lambda_i)$ is the history-dependent conditional crossing probability, which is the chance that, given a first time crossing with $\lambda_i$ since leaving state $A$, $\lambda_{i+1}$ will be crossed before $\lambda_A$. The distribution of first crossing points with $\lambda_i$ since leaving state $A$, is generally not identical to the equilibrium distribution at $\lambda_i$. This aspect includes memory in the expression of Eq. 2 and makes it exact. Milestoning and PPTIS can be viewed as approximate ways to determine the crossing probability by removing or reducing the memory dependence between subsequent interfaces, respectively.

Similarly to the rate $k$ (unit 1/time), the permeability $P$ is a kinetic property (unit length/time). The membrane permeability $P$ is defined as the ratio $J/\Delta c$, where $J$ is the net flux through the membrane when a concentration gradient $\Delta c$ is maintained over the membrane in steady state. Firstly, in the counting method, this ratio is evaluated by counting the membrane crossings in both positive and negative directions and by evaluating the ratio $P = (J^+ + |J^-|)/(2c_{\text{ref}})$ of the bidirectional flux $J^+ + |J^-|$ through the whole membrane and the reference concentration $c_{\text{ref}}$ in a long equilibrium MD simulation. Secondly, in the inhomogeneous solubility-diffusion model, the permeability is estimated assuming diffusive transport, as described by the Smoluchowski equation, giving

$$\frac{1}{P} = e^{-\beta F_{\text{ref}}} \int_{-h/2}^{h/2} \frac{1}{e^{-\beta F(z)} D(z)} \, dz, \qquad (3)$$

with $\beta = 1/(k_B T)$ the inverse temperature, $k_B$ the Boltzmann constant, $T$ the temperature, $h$ the membrane thickness, and $F_{\text{ref}}$ the reference free energy in the solvent phase in region A. The position-dependent free energy $F(z)$ and diffusion $D(z)$ profiles along the membrane normal to that figure in the Smoluchowski equation may be fitted from equilibrium MD with Bayesian analysis (35,36) or a maximum likelihood estimation (37).

Let us now return to the path-sampling methodology. The ratio $P = J^+/c_{\text{ref}}$ is evaluated, where $J^+$ is the flux in the positive direction and $c_{\text{ref}}$ is the reference concentration in state A at the left of the membrane (see Fig. 1). As shown in (17), the definition of $J^+$ has similarities to the RETIS rate $k$,

$$
\begin{aligned}
J^+ &= \frac{\#(A \to M \to B)_{\text{all perm.}}}{T\sigma} \\
k &= \frac{\#(A \to M \to B)_{\text{target}}}{T_{\mathscr{A}}}
\end{aligned}
\tag{4}
$$

where $M$ refers to the membrane region, $T$ is the simulation time of a very long equilibrium simulation, $T_{\mathscr{A}}$ is the part of simulation time spent in overall state $\mathscr{A}$, and $\sigma$ is the cross-section area of the membrane. By juggling Eqs. (4) and (1), it follows that

$$
J^+ = f_A\, p_{\mathscr{A}}\, \frac{N_p}{\sigma}\, P_A(\lambda_B|\lambda_A)
\tag{5}
$$

with $N_p$ the number of permeants in the simulation box and $p_{\mathscr{A}} = T_{\mathscr{A}}/T$. The quantity $p_{\mathscr{A}}$ in Eq. 5 is the probability that the permeants were last in state A rather than in state B. Its evaluation would necessitate full sampling of $\mathscr{B}$, which is, however, not sampled at all in the RETIS simulation! Fortunately, this factor conveniently drops out when evaluating the product $f_A p_{\mathscr{A}}$. This is a key point in the derivation of the practical permeability formula (see (17)). In a last step, the reference concentration $c_{\text{ref}}$ enters in the ratio $J^+/c_{\text{ref}}$. It negates the $N_p/\sigma$ factor in Eq. 5. When combined with the product $f_A p_{\mathscr{A}}$, the reference concentration contribution to the permeability formula becomes a matter of counting the time spent in a user-chosen reference interval $\Delta z$ in state A. With the details given in (17), this gives an equation for the permeability that purely uses RETIS quantities,

$$
P = \frac{\Delta z}{\tau_{\text{ref},[0^-]}} P_A(\lambda_B|\lambda_A)
\tag{6}
$$

where $\tau_{\text{ref},[0^-]}$ is the time spent in $\Delta z$, per path in the $[0^-]$ ensemble.

Yet, Eq. 6 is not straightforward to use with RETIS in practical simulations. The bulk phases at each side of the membrane are, in principle, unbounded such that the order parameter $\lambda$ can have any value from minus infinite to infinite. The application of periodic boundary conditions will prevent this in practice, but could potentially introduce artificial transitions where a permeant ends up at the other side of the membrane without actually moving through it. Ghysels et al. (17) solved both issues by introducing an extra interface $\lambda_{-1} < \lambda_0$ that bounds region A to the left (Fig. 1). Time is frozen for particles that reach beyond $\lambda_{-1}$.

When using the $\lambda_{-1}$ interface, the $[0^-]$ ensemble is replaced with the $[0^{-\prime}]$ path ensemble. Whereas $[0^-]$ only contains paths starting and ending at $\lambda_0$, the $[0^{-\prime}]$ ensemble will also contain paths that start or end at the other side of the A′ region, at $\lambda_{-1}$. Consequently, this changes the number of paths in $[0^{-\prime}]$ vs. $[0^-]$ that could be cut off a long equilibrium trajectory. This in turn affects the time spent *per path* in the $\Delta z$ reference interval by a factor $\xi$, $\tau_{\text{ref},[0^-]}\,\xi = \tau_{\text{ref},[0^{-\prime}]}$ with

$$
\xi = \frac{N_{\to R,[0^{-\prime}]}}{N_{[0^{-\prime}]}}
\tag{7}
$$

The factor $\xi$ expresses this ratio in the number of paths between ensembles, i.e., only the paths arriving to the right at $\lambda_0$ are counted, versus all paths arriving at either $\lambda_{-1}$ or $\lambda_0$ are counted. Using the factor $\xi$, this leads to the final permeability formula in presence of the $\lambda_{-1}$ interface,

$$
P = \frac{\xi \Delta z}{\tau_{\text{ref},[0^{-\prime}]}} P_A(\lambda_B|\lambda_A)
\tag{8}
$$

Alternatively, one can write (17).

$$
P = \frac{\xi P_A(\lambda_B|\lambda_A)}{(\rho_{\text{ref}})_{A'}\, \tau_{[0^{-\prime}]}}
\tag{9}
$$

where $\tau_{[0^{-\prime}]}$ is the average pathlength of paths in ensemble $[0^{-\prime}]$ and $(\rho_{\text{ref}})_{A'}$ is the conditional probability density of the reference region provided that the system is inside state A′. If both $\lambda_0$ and $\lambda_{-1}$ are in the bulk where the free energy is flat, then $(\rho_{\text{ref}})_{A'} = 1/(\lambda_0 - \lambda_{-1})$.

## Improvements of sampling and computational efficiency

Reformulating the permeability expression in RETIS terminology has the obvious advantage that many recent developments in the RETIS method can now be used for permeation simulations. Recently, there have been some interesting advancements in the exact RETIS approach, but even further acceleration while maintaining a good accuracy is possible by reducing the memory dependency of the methodology via a PPTIS-like description of the crossing probabilities. In the next section, we present some simulation results on the combination of replica exchange and PPTIS, coined REPPTIS, in a highly simplified didactical model showing both the importance of replica exchange and memory. In this section, we cover the three aspects by which improvements toward longer timescales are achieved, i.e., development of new MC moves, novel parallelization schemes, and memory reduction.

### MC moves

Like any MC method, the efficiency of the sampling highly depends on the types of moves that are being employed. Until recently, the main MC move in all path-sampling simulations has been the *shooting move* (38) in which a phase point of the previous trajectory is perturbed, usually by a randomization of the velocities alone, after which the equations of motion from this point are integrated forward and backward in time by means of MD until the boundaries of the stable

Vervust et al.

states, $\lambda_A$ or $\lambda_B$, are hit. To ensure detailed balance, the final trajectory is accepted or rejected using a Metropolis-Hastings scheme (21,39).

The *shifting move*, which adds a few steps at the end and removes a few steps at the start of the path or vice versa, was the most frequently executed move in the original TPS method. The standard RETIS (40,41) rate calculation method emerged via TIS (18) from transition path sampling (TPS) (38,42). TIS and RETIS allowed for flexible pathlengths, which made the shifting move both useless and redundant. The *time reversal move*, which simply inverts the direction of time in the old path, used to be employed regularly in path-sampling simulations (TPS, TIS, and RETIS). While it is not so useful in present-day simulation settings where the randomization of velocities in the shooting move is mostly fully randomized, independent of the previous velocities, it is still useful for other types of path analysis such as the predictive capacity identification of reaction triggers (43). RETIS improved the former TPS further via the introduction of the $[0^-]$ path ensemble, and it added the *replica exchange* move between neighboring path ensembles to the palette (41). New advances in path sampling seek to add more alternative moves, replace the main shooting move entirely, or gain greater efficiency by novel parallelization schemes and by optimizing the relative frequency of the moves' execution.

*New MC moves*

In particular, in (17) we added two MC moves, the mirror move and the target swap move, specifically for permeation simulations. The *target swap move* improves the exploration path space whenever more than one permeant is present in the simulation box. As RETIS studies transition from the $\lambda_A$ to the $\lambda_B$ interface defined by the $z$ coordinate of a single target particle, $\lambda = z_t$, the presence of other permeants in the system merely affects the environment of the target permeant. The target swap move, however, uses the statistics of the other permeants more effectively by a random reassignment of the target. This new target permeant might be located in a different area of the simulation box, and thus higher sampling decorrelation is likely achieved. The *mirror move* (17) increases the sampling in periodic systems by completely mirroring the particle's coordinates in the $xy$ plane, effectively changing the reaction coordinate from $\lambda = z_t$ to $\lambda = -z_t$ with respect to the original coordinates. This implies that permeation pathways through the membrane in both directions are sampled, which also improves sampling efficiency. Despite the fact that the target swap move and the mirror move only operate in the $[0^{-\prime}]$ path ensemble, the faster exploration in the directions orthogonal to the reaction coordinate are felt by all the other path ensembles due to the replica exchange moves between path ensembles, as was clearly illustrated in a membrane model with two unequal permeation channels (17).

Another promising trend is to change the main MC move itself via so-called *subtrajectory moves* (30,31). The main idea

behind these moves is that successive paths, created by the shooting move, are correlated, which leads to a statistical inefficiency, $\mathcal{N}$, of several tens or hundreds of paths (40). Not saving every path, but saving every $N_s$-th path, with $N_s < \mathcal{N}$, will typically not lead to any loss in statistical precision in the final result as a lower number of stored paths, which is used in the analysis, is compensated by a reduction in the correlation between the paths that are saved. While this may be a good strategy to save disk capacity and time required for writing to disk, the subtrajectory moves go a step further by significantly reducing the number of MD steps for paths that do not need to be saved. The subtrajectory move is best combined with the *high-acceptance* technique (30,31).

*Parallellization*

Further efficiency gains without invoking any approximation or Markovian assumption can be achieved via a smart parallelization scheme and by maximizing the replica exchange swapping frequency (32). Parallel computing will typically distribute the same number of processing units per ensemble to carry out the computational intensive standard moves. This makes the parallelization of the RETIS method a nontrivial task as each path can have a different length and the average pathlength differs for each path ensemble. Standard RETIS simulations apply the replica exchange swapping moves and standard MC moves alternately. The swapping move is cheap, but requires that the ensembles involved in the swap have completed their previous move. This means that, if the standard moves in each ensemble require different computing times, several processing units have to wait for the slow ones to finish, i.e., the replicas are cost-unbalanced. Roet et al. (32) solve this problem through a fundamentally new approach to the generic replica exchange method in which ensembles are not updated in cohort.

Roet et al. (32) also show that the number of replica exchange moves, in between two shooting or subtrajectory moves, can effectively be set to infinite without having to do an infinite number replica exchange moves explicitly. While the idea of infinite swapping has been suggested before (44–47), a reformulation of the implicit infinite swapping problem in terms of permanents allows for a much better scaling with the number of interfaces. The non-cohort infinite replica exchange approach applied to RETIS, coined $\infty$RETIS, opens the way for massively parallel path-sampling simulations for computing rate constants (40), activation energies (48), permeability constants (17), and mechanistic analysis for reaction triggers (43,49).

*Reduction in memory*

Still, when the individual trajectories themselves are too long to be simulated, the statistics of long trajectories should be obtained via shorter ones without actually sampling any trajectory going all the way from $\lambda_A$ to $\lambda_B$. This is essentially the idea behind milestoning (19) and PPTIS (23). This strategy will generally cause the method to be no longer exact

unless the interfaces are isocommittor surfaces (24). However, the isocommittor surfaces are generally not known and extremely difficult and costly to determine via simulations. The lack of knowledge about the isocommittor can be compensated by adding a bit of memory to the interface crossing probabilities. We denote the PPTIS path ensembles as $[i^\pm]$ (23). Trajectories in path ensemble $[i^\pm]$ are restricted by the $\lambda_{i-1}$ and $\lambda_{i+1}$ interfaces. They can start and end at either side, but should at least cross the middle interface $\lambda_i$ once. From these path ensembles, two-directional *local* crossing probabilities are obtained, $p_i^\pm$, $p_i^=$, $p_i^\mp$, and $p_i^\ddagger$. Here, the lower sign refers to the past conditional direction and the upper sign refers to the measure of the probability in the future, measured from a point in time where $\lambda_i$ is crossed for the first time since its latest crossing with either $\lambda_{i-1}$ or $\lambda_{i+1}$. For instance, both $p_i^\pm$ and $p_i^\ddagger$ refer to the probability that $\lambda_{i+1}$ is crossed earlier than $\lambda_{i-1}$ (future condition) after a first crossing with $\lambda_i$. But their past condition is different and equal to "given it came from $\lambda_{i-1}$" and "given it came from $\lambda_{i+1}$," respectively. The local crossing probabilities with the same past condition add up to one: $p_i^\pm + p^= = p_i^\mp + p_i^\ddagger = 1$. Once a sufficient number of paths in the $[i^\pm]$ ensemble is sampled, the local crossing probability is determined by simply counting the appropriate paths with specific future and past conditions, e.g., $p_i^\pm$ is given by the number of paths starting at $\lambda_{i-1}$ *and* ending at $\lambda_{i+1}$ divided by the number of paths starting at $\lambda_{i-1}$.

The PPTIS formalism is based on recursive relations where the local crossing probabilities are linked to *global* crossing probabilities:

$$P_j^+ \approx \frac{p_{j-1}^\pm P_{j-1}^+}{p_{j-1}^\pm + p_{j-1}^= P_{j-1}^+}, P_j^- \approx \frac{p_{j-1}^\mp P_{j-1}^-}{p_{j-1}^\pm + p_{j-1}^= P_{j-1}^-} \quad (10)$$

$$P_1^+ = P_1^- = 1$$

where $P_j^+ = P_A(\lambda_j | \lambda_1)$ is the chance to cross $\lambda_j$ before $\lambda_A = \lambda_0$ given that $\lambda_1$ is crossed at this moment while $\lambda_A$ was crossed more recently than $\lambda_1$. Similarly, $P_j^-$ is the chance that $\lambda_A$ is crossed before $\lambda_j$ given that $\lambda_{j-1}$ is crossed at this moment while $\lambda_j$ was crossed more recently than $\lambda_{j-1}$. From these recursive relations, the *overall* crossing probability from $\lambda_A$ to $\lambda_B$ can be computed from

$$P_A(\lambda_B | \lambda_A) = P_A(\lambda_n | \lambda_0) = P_A(\lambda_1 | \lambda_0) P_A(\lambda_n | \lambda_1) = p_0^\pm P_n^+ \quad (11)$$

Here, the $p_0^\pm$ probability is slightly different from the $p_i^\pm$ definitions with $i > 0$ in the sense that it is just the probability to reach $\lambda_1$ before $\lambda_0$ after a positive crossing with $\lambda_0$, i.e., $p_0^\pm$ has no additional past condition.

The larger the distance between interfaces, the more memory is included in the calculation, the more accurate is Eq. (10). The calculation of memory loss functions is a way to estimate the required distance between interfaces (23). On the other hand, the interfaces should be placed relatively close to each other to obtain the best efficiency. This can lead to conflicting strategies for parameter optimization. A potential solution could be to use path history beyond the boundaries of the $[i^\pm]$ ensemble. This kind of information could in principle become available if PPTIS is also combined with replica exchange moves.

The potential of performing replica exchange between path ensembles was first suggested for PPTIS (48). Yet, this idea has so far never been put in practice. The replica exchange move $[i^\pm] \leftrightarrow [(i+1)^\pm]$ in PPTIS is more costly than the swapping move $[i^+] \leftrightarrow [(i+1)^+]$ in RETIS. In RETIS, full trajectories are swapped without the need to do additional MD steps. In PPTIS, it is first checked whether the $[i^\pm]$ path ends at $\lambda_{i+1}$ and whether the $[(i+1)^\pm]$ path starts at $\lambda_i$. If not, the move is directly rejected. However, if so, the $[i^\pm]$ and $[(i+1)^\pm]$ paths are extended forward and backward in time, respectively. Subsequently, the extended paths are trimmed in accordance to the new path ensemble boundaries to which the paths are being transferred to.

This article presents the first applications of the replica exchange and PPTIS combination, which we coin REPPTIS. The algorithms are implemented in the PyRETIS code, which is readily available to be used in combination with other MD simulation packages such as GROMACS, OpenMM, or CP2K (50,51). The first application is a model system, showing both the importance of the replica exchange moves and the effect of memory. In the second application of ibuprofen permeation, we show how RETIS is challenged by metastable states, which can make the paths prohibitively long, whereas REPPTIS can be used to simulate full membrane transits. The results of these two examples are presented in the next section.

## RESULTS

### Permeation through a maze potential

A two-dimensional toy system is developed to demonstrate the role of memory in permeability calculations. A Langevin particle is permeating along the $z$ direction from the water phase through the membrane (Fig. 2). The propagation of the permeation is measured by the $z$ coordinate of the particle, so the order parameter is $\lambda = z$. The coordinate $y$ describes the orthogonal degrees of freedom, which could be a general coordinate such as the orientation of the molecule or the local composition of heterogeneous membranes. Here, a membrane is chosen with different permeation pathways; for instance, corresponding to different regions of a heterogeneous membrane. The membrane is represented by a maze potential with two permeation channels (upper channel for $y > 0.5$, lower channel for $y < 0.5$). Passage through the lower channel is entropically unfavorable, as the lower channel is only accessible via an aperture at about $z = 0.44$. Passage through the upper channel requires the Langevin particle to
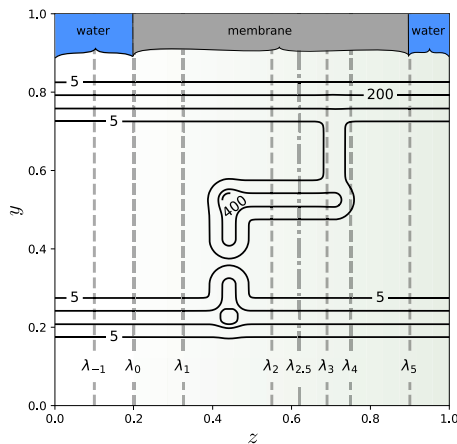
Vervust et al.



FIGURE 2 The potential energy $V(z,y)$ for the two-dimensional maze system. The isolines display the specific values for the potential. In addition, a green gradient is added in the horizontal direction from $z = 0.2$ to $z = 1.0$ to visualize the small, positive, and linear tilt with a slope of 0.5 that exists within that range in the potential. Interfaces $\lambda_{-1}, \dots, \lambda_5$ are indicated with vertical dashed lines. Reduced units are used. To see this figure in color, go online.

overcome an energy barrier at $z = 0.71$. The maze can be thought of as a pinball machine, where the ball can either go through the "flippers" of the lower channel, or go over the "bump" in the upper channel. When the initial path is located in the upper channel passing over the bump, the orthogonal degree of freedom ($y$) will need to be sampled sufficiently with path sampling to detect the alternative pathway through the flippers, and vice versa. Moreover, for every $z$ value, there is a broad $y$ region in the upper and/or lower channel where the ball can easily move locally somewhat left or right. This local picture could give the impression that the ball has no entropic nor energetic challenges to overcome at all. However, all *complete* permeation pathways will need to overcome the entropic or energetic barrier, so including all memory in the pathways gives a distinctly different picture than the local picture. This makes our test case a good illustration of the role of memory. If one focuses on pieces of trajectories that only move somewhat to the left or right, i.e., if memory is too short, one will miss the dynamics included in the complete pathways.

We do different types of simulations.

- RETIS simulation, which retains all memory and thus gives exact kinetics. This is the benchmark.
- PPTIS simulation, where memory is reduced. In the $[i^\pm]$ ensemble, paths that cross $\lambda_i$ are cut short when they pass a neighboring interface at $\lambda_{i-1}$ or $\lambda_{i+1}$.

- REPPTIS simulation, where memory is reduced and where replica exchange moves between the $[i^\pm]$ ensembles are allowed. This potentially incorporates additional memory (see previous section).

Six interfaces were chosen along the $z$ axis ($\lambda_A = \lambda_0, \dots, \lambda_5 = \lambda_B$) in the membrane region between $\lambda_A = 0.2$ and $\lambda_B = 0.9$. The $[0^+], \dots, [4^+]$ path ensembles were sampled for a total of 100,000 MC moves with the PyRETIS code (50,51) using shooting moves, the wire fencing (31) variant of the subtrajectory moves with 6 subpaths, and replica exchange moves. The $[i^\pm]$ ensembles were sampled for (RE)PPTIS without wire fencing moves. In addition, the $\lambda_{-1}$ interface was used at $z = 0.1$ to bound the region at the left of the membrane. The properties $\xi$ and $\tau_{\text{ref},[0^{-\prime}]}$ with the reference interval set to [0.1,0.2] were computed from the $[0^{-\prime}]$ path ensemble. Two PPTIS simulations were run, where the first was initialized with a reactive path through the lower entropic barrier, and the second with a reactive path through the upper energetic barrier. These PPTIS simulations are referred to as PPTIS 1 and PPTIS 2, respectively.

To challenge the memory reduction, we also run REPPTIS with an extra interface at $z = 0.62$ between $\lambda_2$ and $\lambda_3$, which we refer to as $\lambda_{2.5}$ in this text for convenience. For RETIS, extra interfaces typically increase the accuracy as more paths are sampled, and matching the probabilities can be done with higher accuracy. For (RE)PPTIS, however, the extra interface implies a more drastic cut in memory as some of the $[i^\pm]$ ensembles will span a smaller spatial area. The simulation without and with the extra $\lambda_{2.5}$ interface are referred to as REPPTIS 1 and REPPTIS 2, respectively.

A tilt potential with slope 0.5 was superimposed on the maze potential (green gradient in Fig. 2) for $z \geq 0.2$, mimicking a membrane barrier. Details about setup, simulations, and code to generate more general maze potentials are given in the supporting material. Reduced units are used, and reported errors are standard errors based on block averaging.

## The maze: Effect of memory and replica exchange move

The average pathlength in the highest RETIS ensemble $[4^+]$ is 47.5, while the $[4^\pm]$ (RE)PPTIS ensembles have an average pathlength of about 3.8. The goal of reducing the length of the MD trajectories by memory reduction is thus clearly achieved. With Eq. 8, the permeability from RETIS is equal to $2.54 \times 10^{-5}$ ($\pm 8\%$). To compare the effect of memory on this permeability value, this discussion will focus on the crossing probability $P_A(\lambda_B | \lambda_A)$, which is the only factor in Eq. 8 that may be affected by the memory reduction. The overall crossing probability is given in Table 1. We first discuss the simulations without the extra $\lambda_{2.5}$ interface. The PPTIS 1 and PPTIS 2 simulations

**TABLE 1** Overall crossing probability through the membrane from $\lambda_A = 0.2$ to $\lambda_B = 0.9$, without and with extra interface $\lambda_{2.5}$

| Simulation | Overall $P_A(\lambda_B|\lambda_A)$ [$10^{-4}$] | Forward | | | Backward | | |
|---|---|---|---|---|---|---|---|
| | | $p_2^\pm$ | $p_{2.5}^\pm$ | $p_3^\pm$ | $p_2^\mp$ | $p_{2.5}^\mp$ | $p_3^\mp$ |
| RETIS | 2.65 ($\pm 5\%$) | | | | | | |
| PPTIS 1 initial lower | 6.44 ($\pm 17\%$) | 0.41 | | 0.54 | 0.45 | | 0.63 |
| PPTIS 2 initial upper | 0.93 ($\pm 34\%$) | 0.18 | | 0.24 | 0.70 | | 0.96 |
| REPPTIS 1 | 2.14 ($\pm 12\%$) | 0.19 | | 0.47 | 0.56 | | 0.66 |
| REPPTIS 2 with $\lambda_{2.5}$ | 2.94 ($\pm 9\%$) | 0.59 | 0.28 | 0.43 | 0.39 | 0.79 | 0.75 |

For PPTIS, the initial path was either in the upper (PPTIS 1) or lower (PPTIS 2) channel. The local crossing probabilities $p^\pm$ and $p^\mp$ are given for some of the $[i^\pm]$ ensembles.

significantly overestimate and underestimate the crossing probability, respectively, by more than a factor of 2. The added replica exchange moves in REPPTIS 1 improve the crossing probability considerably. Fig. 3 plots the intermediate crossing probabilities $P_A(\lambda_i|\lambda_A)$ to reach $\lambda_i$, showing that the first deviations between the simulations start at $\lambda_3$, when the Langevin particle has entered the maze.

Let us look at the origin of these deviations by tracing some randomly selected exemplary paths in the different ensembles in Fig. 4. In the reference simulation, RETIS, the reactive paths in $[4^+]$ cross both the entropic and energetic barrier, where the particle prefers the "flippers" channel (lower) rather than the "bump" channel (upper).

In PPTIS 1, the initial path is located in the lower channel and it remains there indefinitely in the $[3^\pm]$ ensemble. The MC shooting moves in $[3^\pm]$ cannot result in a switch to the other channel, so the path is stuck. The absence of such nonlocal moves that allow channel switching breaks the ergodicity of the sampling. The effect of this sampling deficit is modest as can be concluded from the $p_3^\pm$ and $p_3^\mp$
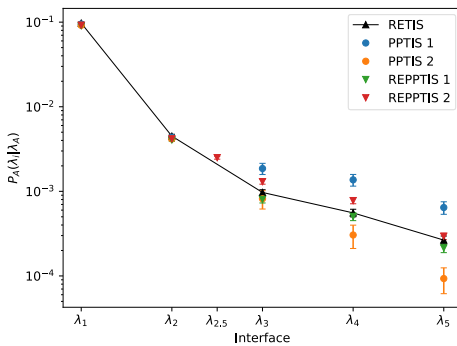


**FIGURE 3** Crossing probabilities $P_A(\lambda_i|\lambda_A)$ for RETIS (*black*), PPTIS (*circles*), and REPPTIS (*triangles*) simulations. The RETIS simulation provides the benchmark values, which are connected by a line. REPPTIS 2 includes the extra $\lambda_{2.5}$ interface. Values to the right are the global crossing probabilities $P_A(\lambda_5|\lambda_A) = P_A(\lambda_B|\lambda_A)$. The error bars are standard errors based on block averaging. To see this figure in color, go online.

values in Table 1 since crossings with $\lambda_3$ are more likely in the lower channel where the potential energy is low. Also in the $[2^\pm]$ ensemble, the paths remain for most of the simulation in the lower channel, but here the ergodicity problem is more severe since most first time crossings with $\lambda_2$ coming from $\lambda_1$ should be in the upper channel. As it is more difficult to reach $\lambda_3$ from $\lambda_2$ in the undersampled upper channel, $p_2^\pm$ is overestimated. For PPTIS 2, the initial path is in the upper channel and the absence of nonlocal MC moves again breaks the ergodicity of the sampling, where the paths in $[3^\pm]$ are now stuck in the upper channel. In REPPTIS 1, we have added the replica exchange moves and, impressively, this move reintroduces ergodicity. When the $y$ axis is sampled in the $[0^{-\prime}]$ ensemble, this effect can be transported to the other ensembles with the swap move. Adding the exchange move will thus effectively allow for switching between channels.

The crossing probabilities of PPTIS are affected by the particle being stuck in $[3^\pm]$ (or somewhat stuck in $[2^\pm]$) in a particular channel. It is much easier to reach $\lambda_3$ in the lower channel instead of the upper channel, where $\lambda_3$ is high uphill on the wall's slope. The true mechanism has a mixture of both pathways (see RETIS). This gives an overestimation of the $\lambda_1 \rightarrow \lambda_2 \rightarrow \lambda_3$ and $\lambda_2 \rightarrow \lambda_3 \rightarrow \lambda_4$ crossing probability by PPTIS 1 and an underestimation by PPTIS 2. Numerically, this is also reflected in the local crossing probabilities. A selection is shown in Table 1; other local crossing probabilities were not statistically different between the simulations, as expected. PPTIS 1 is mainly located in the easier flipper channel, and $p_2^\pm$ is overestimated by PPTIS 1 compared with REPPTIS 1 (0.41 vs. 0.19), which increases the overall crossing probability in PPTIS 1. Likewise, PPTIS 2 is mainly located in the more difficult bump channel, and $p_3^\pm$ is strongly underestimated by PPTIS 2 compared with REPPTIS 1 (0.24 vs. 0.47). In combination with the higher backwards $p_2^\mp$ and $p_3^\mp$ values, this results in a lower overall crossing probability. In a very long PPTIS simulation, the $[2^\pm]$ ensemble could eventually be correctly sampled, but the $[3^\pm]$ ensembles will remain stuck with paths resembling the initial path.

Finally, we discuss REPPTIS 2 with the extra interface at $\lambda_{2.5}$, which further reduces memory. Fig. 3 shows that the probability to reach $\lambda_3$ is overestimated, which can be expected because of tunneling between $[2^\pm]$ and $[2.5^\pm]$. In $[2^\pm]$, the particle can move freely in the upper channel, as it is not hindered by the energetic barrier located to the right of $\lambda_{2.5}$. In $[2.5^\pm]$, the particle can move freely in the lower channel, as it is not hindered by the entropic barrier located to the left of $\lambda_2$. This can also be seen in Fig. 4, where the LMR and RML example paths of $[2^\pm]$ and $[2.5^\pm]$ are predominantly located in the upper and lower channel, respectively. Connecting these two ensembles to derive a $\lambda_1 \rightarrow \lambda_2 \rightarrow \lambda_3$ crossing probability, the particle seems to switch from the upper channel in $[2^\pm]$ to the lower channel
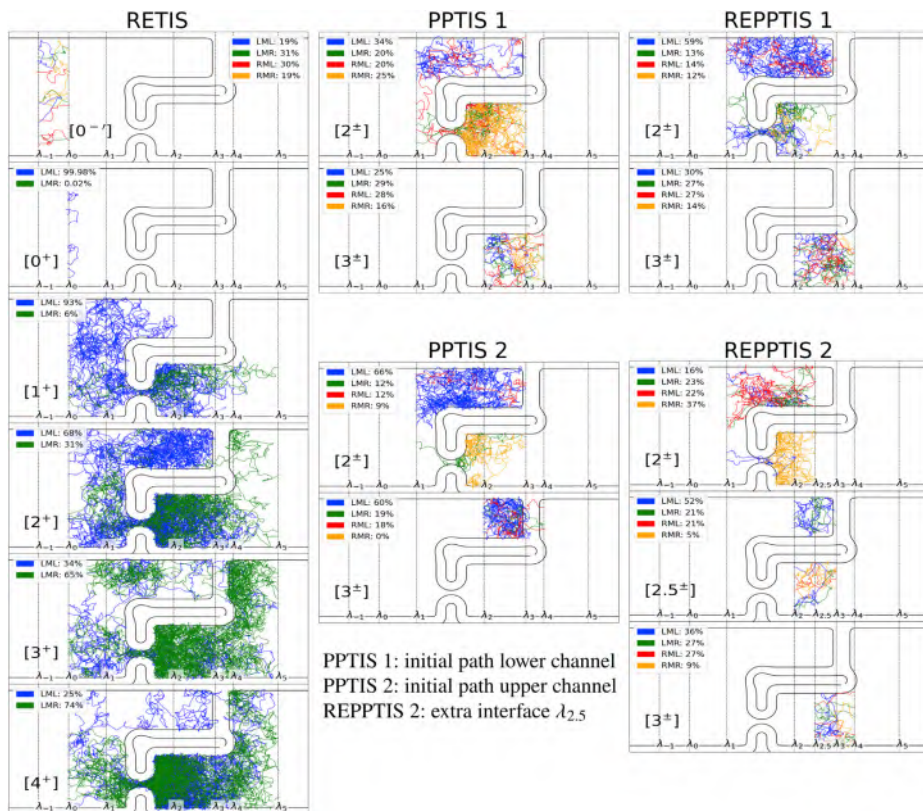
Vervust et al.



FIGURE 4  Example paths of each simulation. The paths are colored according to the four possible path types, which are determined by the interfaces where the path begins and ends. For example, if a path of the $[i^{\pm}]$ starts at interface $\lambda_{i-1}$ and ends at interface $\lambda_{i+1}$, it is labeled an LMR path as it started from the left (L) of $\lambda_i$ and ended at the right (R) of $\lambda_i$. The four possibilities are LML paths (*blue*), LMR paths (*green*), RML paths (*red*), and RMR paths (*orange*). The weight of each path type is indicated in percentages. For each ensemble, 10 paths are randomly selected, respecting the weight of the path types. For example, of the 10 paths in the $[2^{\pm}]$ PPTIS 1 ensemble, 3 are LML (34%), 2 are LMR (20%), 2 are RML (20%), and 3 are RMR (25%). To see this figure in color, go online.

in $[3^{\pm}]$, as if it had tunneled through the wall. In other words, the particle "forgets" its passage through the flippers. A too harsh reduction of memory can thus lead to an overestimate of the permeability. For REPPTIS 1, no such tunneling occurs between $[2^{\pm}]$ and $[3^{\pm}]$, because $\lambda_3$ is located on the rising edge of the energetic barrier. Surprisingly, the REPPTIS 2 crossing probability decreases at the end, resulting in a total crossing probability that lies close to RETIS. This is, however, a lucky cancellation of errors as tunneling also happens from right to left. The paths of RETIS and REPPTIS 1 have a small probability of recrossing the energetic or entropic barriers, while the reverse

tunneling in REPPTIS 2 makes this more likely, which results in a decrease of the global crossing probability.

### Permeation of ibuprofen drug molecule

Whereas the maze system was built to showcase the role of memory, we now study an application that is more representative for a typical membrane permeability simulation. Passive permeability through (lipid) membranes is of vital importance for drug design, as it gives insight to the timescale at which drugs transit the membrane for a given concentration gradient (6,52). The nonsteroidal anti-inflammatory

120

drug ibuprofen has to cross several membranes before realizing its inhibitory effect on the cyclooxygenase enzymes COX-1 and COX-2 (53,54). Path sampling is now used to investigate the permeation of ibuprofen through a dioleoyl-phosphatidylcholine (DOPC) membrane. The presence of metastable states combined with an orthogonal degree of freedom will put REPPTIS to the test.

Assume $z$ is the center-of-mass distance (in unit nanometers) of ibuprofen to the bilayer midplane, and $\theta$ is the dihedral angle determining the OH orientation in the carboxyl group of ibuprofen. The free energy profile $F(z,\theta)$ of ibuprofen in a DOPC bilayer is shown in Fig. 5 A, which was recreated from data in (55). Details of the implementation can be found in the supporting material. The two stable OH bond conformations are *cis* ($\theta \approx$ 0) and *trans* ($\theta \approx \pi$), which are visualized in Fig. 5 B. The
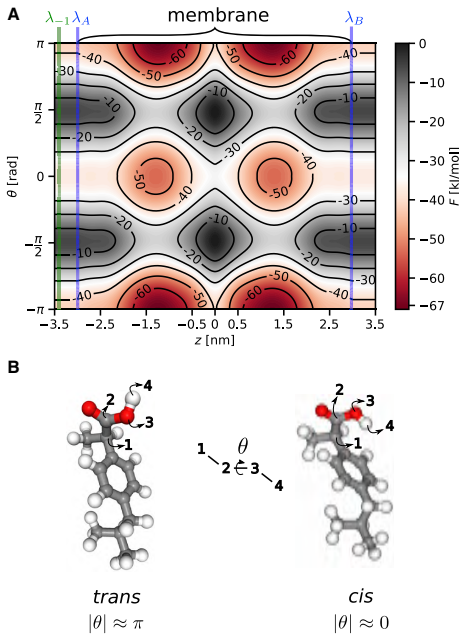


**A**

**B**

*trans*
$|\theta| \approx \pi$

*cis*
$|\theta| \approx 0$

FIGURE 5 (*A*) Two-dimensional free energy profile $F(z,\theta)$ of ibuprofen in a DOPC bilayer. $F$ is periodic in $\theta$ with period $2\pi$. The regions $\theta \approx 0$ and $|\theta| \approx \pi$ are the *cis* and *trans* configurations, respectively. The membrane is located in $z \in [-3, 3]$, while $z < -3$ and $z > 3$ represent the water phases near each leaflet. The RETIS simulation to calculate $P_{entr}$ uses the $\lambda_{-1}$ interface (*green vertical line*). The simulation domains for all other RE(PP)TIS simulations lie in between the $\lambda_A$ and $\lambda_B$ interfaces (*blue lines*) of the full permeation REPPTIS simulation. (*B*) Ibuprofen in the *cis* and *trans* configurations. The atoms that define the dihedral angle $\theta$ of the hydroxyl hydrogen of the carboxyl group are annotated with arrows. To see this figure in color, go online.

free energy profile contains minima in each leaflet ($z \in [-1.5, -1]$ and $z \in [1, 1.5]$), such that a full transit across the bilayer consists of at least two successive rare events. This means that RETIS is no viable option to study the full permeation event, as the paths would get stuck in the metastable states, resulting in extremely long paths. The memory reduction in REPPTIS greatly reduces the pathlengths and enables the study of the full permeation using a single simulation.

Permeation of ibuprofen through the DOPC bilayer is characterized by three steps, i.e., entering the membrane, hopping over an internal barrier, and escaping from the membrane. Starting from the water phase in $z < -3$, ibuprofen *enters* the energy minimum around $z \in [-1.5, -1.0]$ in the first leaflet. Given that ibuprofen crosses $z = -3$ to the right, the molecule has a crossing probability denoted by $P_{entr}$ to enter into this stable region. This probability is not 1, as friction by the membrane molecules can let the molecule return to $-3$, rather than fully entering. Next, ibuprofen must overcome the *internal* barrier between the leaflets to reach the second energy minimum $z \in [1, 1.5]$. The corresponding crossing probability over the internal barrier is denoted $P_{int}$. Finally, ibuprofen needs to *escape* from the second stable region and reach the water phase ($z > 3$). This crossing probability is denoted $P_{esc}$.

These characteristic crossing probabilities $P_{entr}$, $P_{int}$, and $P_{esc}$ are calculated using RETIS and REPPTIS (simulation details in supporting material), and are given in Table 2. Each simulation was run twice, where the initial path was either in the *cis* or *trans* configuration. As it was verified that transitions between these configurations happened in all of the ensembles, the data of both runs were merged.

As the paths would become too long to simulate a full transit with RETIS, the three characteristic crossing probabilities were used in a Markov model to estimate the full transit probability $P_{trans}$ from $z = -3$ to $z = 3$. Let $k_{int} = f_{int} P_{int}$ and $k_{esc} = f_{esc} P_{esc}$ be the internal and escape rates, respectively, where the fluxes $f_{int}$ and $f_{esc}$ are part of the RETIS simulation output. As shown in the supporting material, the transit probability is approximated by $P_{trans} \approx (P_{entr} k_{int})/(k_{esc} + 2k_{int})$. In contrast to RETIS, REPPTIS is capable of calculating $P_{trans}$ of the full membrane transit using a single simulation. Both the REPPTIS and the approximate Markov RETIS values of $P_{trans}$ are given in Table 2.

**TABLE 2 Characteristic crossing probabilities of ibuprofen permeation through a phospholipid bilayer**

| Simulation | $P_{entr}$ [$10^{-2}$] | $P_{int}$ [$10^{-6}$] | $P_{esc}$ [$10^{-6}$] | $P_{trans}$ [$10^{-3}$] |
|---|---|---|---|---|
| RETIS | 17 ($\pm 7\%$) | 1.2 ($\pm 11\%$) | 2.3 ($\pm 8\%$) | 5.4 ($\pm 23\%$) |
| REPPTIS | | 1.1 ($\pm 8\%$) | 2.2 ($\pm 5\%$) | 6.1 ($\pm 7\%$) |

$P_{trans}$ is the crossing probability of a full membrane transit: the REPPTIS value is from a simulation, while the RETIS value is an estimate based on a simple Markov model. Reported errors are standard errors based on block averaging.

Vervust et al.

The RETIS and REPPTIS simulations result in statistically equivalent crossing probabilities for both the internal and escape transitions. From the entrance RETIS simulation, the factor $(\xi \Delta z)/\tau_{\mathrm{ref},[0^{-\prime}]} = (5.0 \pm 1\%) \times 10^{-2}$ nm/ps is obtained, which enters the permeability equation (Eq. 8). Using the Markov model with the characteristic crossing probabilities of the RETIS simulations, the permeability of ibuprofen becomes $(27 \pm 23\%)$ cm/s. Using $P_{\mathrm{trans}}$ of the full permeation REPPTIS simulation, the permeability of ibuprofen is estimated to be $(30 \pm 7\%)$ cm/s. This value, based on the two-dimensional $F(z, \theta)$ profile at 303 K of (55), is in reasonable agreement with the ibuprofen permeability $(92 \pm 6)$ cm/s through dipalmitoylphosphatidylcholine at 323 K as obtained from MD simulations and the inhomogeneous solubility-diffusion model (56).

## CONCLUSION

In this article we first reviewed the recently developed theoretical framework for calculating permeability coefficients using the RETIS methodology. The approach requires a slight modification of the $[0^-]$ path ensemble to the $[0^{-\prime}]$ ensemble, which describes the paths at the left side of $\lambda_A$. The RETIS-based permeability can be computed with exponential reduction in time compared with standard MD, while it still gives exactly the same result without introducing any approximation. The mathematical formulation of microscopic permeability in terms of RETIS properties has the advantage that recent algorithmic developments in the RETIS method can directly be applied, such as the recently developed MC moves for generating new paths more efficiently (17,30,31) and the nonsynchronous replica exchange approach (32).

However, if the individual transition paths themselves are long, it may be wise to give up some of the method's exactness for the sake of obtaining shorter paths. This idea underlies the PPTIS method in which the statistics of long transition paths is obtained via paths with much shorter range using a memory loss assumption. Still, some memory is retained in the conditional local crossing probabilities that are computed. In this article, we combined the PPTIS method with replica exchange into a new implementation, coined the REPPTIS method. We applied PPTIS and REPPTIS on a didactic model and on the permeation of ibuprofen based on a realistic free energy surface. The results showed the importance of both replica exchange and memory, as simulations without them gave wrong permeability estimates.

There are several interesting opportunities to improve the REPPTIS method further. Note that the extension of paths by means of MD in a swapping move could yield additional information that adds back in some of the lost memory. Before the extended trajectories are trimmed to fit the boundaries of the new path ensemble, these extensions provide continuous trajectories that go beyond the range of three consecutive interfaces. Using the information of untrimmed trajectories might be exploited in future variants of the REPPTIS method since it could solve the conflicting benefits of having interfaces close enough for efficiency and far apart for accuracy. Moreover, the time information is not yet exploited by TIS-based methods. For instance, the PPTIS crossing probabilities relate to the chance that a specific interface is crossed before another irrespective how long it takes. Another future development that we want to achieve is the inclusion of time durations in the statistical description of the crossing probabilities, as is done in milestoning, to compute (conditional) mean first passage times and diffusion coefficients. Note that other milestone variations such as the use of multidimensional interface networks via, e.g., Voronoi cells (57), can in principle be applied within a REPPTIS framework as well. We can therefore conclude that REPPTIS is a promising method to enable permeation simulations with high efficiency and accuracy that might not be easily achieved by any other method.

## SUPPORTING MATERIAL

## AUTHOR CONTRIBUTIONS

## ACKNOWLEDGMENTS

## DECLARATION OF INTERESTS

The authors declare no competing interests.

## REFERENCES

1. Yang, N. J., and M. J. Hinner. 2015. Getting across the cell membrane: an overview for small molecules, peptides, and proteins. *In* Site-Specific Protein Labeling: Methods and Protocols. A. Gautier and M. J. Hinner, eds Springer, pp. 29–53.

2. Awoonor-Williams, E., and C. N. Rowley. 2016. Molecular simulation of nonfacilitated membrane permeation. *Biochim. Biophys. Acta.* 1858:1672–1687.

3. Shinoda, W. 2016. Permeability across lipid membranes. *Biochim. Biophys. Acta.* 1858:2254–2265.

4. Bennion, B. J., N. A. Be, …, T. S. Carpenter. 2017. Predicting a drug's membrane permeability: a computational model validated with in vitro permeability assay data. *J. Phys. Chem. B.* 121:5228–5237.

5. Hannesschlaeger, C., A. Horner, and P. Pohl. 2019. Intrinsic membrane permeability to small molecules. *Chem. Rev.* 119:5922–5953.

6. Menichetti, R., K. H. Kanekal, and T. Bereau. 2019. Drug-membrane permeability across chemical space. *ACS Cent. Sci.* 5:290–298.

7. Levental, I., and E. Lyman. 2023. Regulation of membrane protein structure and function by their lipid nano-environment. *Nat. Rev. Mol. Cell Biol.* 24:79.

8. Dotson, R. J., C. R. Smith, …, S. C. Pias. 2017. Influence of cholesterol on the oxygen permeability of membranes: insight from atomistic simulations. *Biophys. J.* 112:2336–2347.

9. Venable, R. M., A. Krämer, and R. W. Pastor. 2019. Molecular dynamics simulations of membrane permeability. *Chem. Rev.* 119: 5954–5997.

10. Davoudi, S., and A. Ghysels. 2021. Sampling efficiency of the counting method for permeability calculations estimated with the inhomogeneous solubility–diffusion model. *J. Chem. Phys.* 154, 054106.

11. Marrink, S. J., and H. J. C. Berendsen. 1994. Simulation of water transport through a lipid membrane. *J. Phys. Chem.* 98:4155–4168.

12. De Vos, O., R. M. Venable, …, A. Ghysels. 2018. Membrane permeability: characteristic times and lengths for oxygen and a simulation-based test of the inhomogeneous solubility-diffusion model. *J. Chem. Theory Comput.* 14:3811–3824.

13. Ghysels, A., A. Krämer, …, R. W. Pastor. 2019. Permeability of membranes in the liquid ordered and liquid disordered phases. *Nat. Commun.* 10:5616.

14. Torrie, G., and J. Valleau. 1977. Nonphysical sampling distributions in Monte Carlo free-energy estimation: umbrella sampling. *J. Comput. Phys.* 23:187–199.

15. Darve, E., and A. Pohorille. 2001. Calculating free energies using average force. *J. Chem. Phys.* 115:9169–9183.

16. Comer, J., J. C. Gumbart, …, C. Chipot. 2015. The adaptive biasing force method: everything you always wanted to know but were afraid to ask. *J. Phys. Chem. B.* 119:1129–1151.

17. Ghysels, A., S. Roet, …, T. S. van Erp. 2021. Exact non-Markovian permeability from rare event simulations. *Phys. Rev. Res.* 3, 033068.

18. van Erp, T. S., D. Moroni, and P. G. Bolhuis. 2003. A novel path sampling method for the calculation of rate constants. *J. Chem. Phys.* 118:7762–7774.

19. Faradjian, A. K., and R. Elber. 2004. Computing time scales from reaction coordinates by milestoning. *J. Chem. Phys.* 120:10880–10889.

20. Allen, R. J., P. B. Warren, and P. R. ten Wolde. 2005. Sampling rare switching events in biochemical networks. *Phys. Rev. Lett.* 94, 018104.

21. Metropolis, N., A. W. Rosenbluth, …, E. Teller. 1953. Equation of state calculations by fast computing machines. *J. Chem. Phys.* 21: 1087–1092.

22. van Erp, T. 2012. Dynamical rare event simulation techniques for equilibrium and nonequilibrium systems. *Adv. Chem. Phys.* 151:27.

23. Moroni, D., P. G. Bolhuis, and T. S. van Erp. 2004. Rate constant for diffusive processes by partial path sampling. *J. Chem. Phys.* 120:4055–4065.

24. Vanden-Eijnden, E., M. Venturoli, …, R. Elber. 2008. On the assumptions underlying milestoning. *J. Chem. Phys.* 129, 174102.

25. Cardenas, A. E., and R. Elber. 2013. Computational study of peptide permeation through membrane: searching for hidden slow variables. *Mol. Phys.* 111:3565–3578.

26. Cardenas, A. E., and R. Elber. 2014. Modeling kinetics and equilibrium of membranes with fields: milestoning analysis and implication to permeation. *J. Chem. Phys.* 141, 054101.

27. Fathizadeh, A., and R. Elber. 2019. Ion permeation through a phospholipid membrane: transition state, path splitting, and calculation of permeability. *J. Chem. Theory Comput.* 15:720–730.

28. Votapka, L. W., C. T. Lee, and R. E. Amaro. 2016. Two relations to estimate membrane permeability using milestoning. *J. Phys. Chem. B.* 120:8606–8616.

29. Riccardi, E., A. Krämer, …, A. Ghysels. 2021. Permeation rates of oxygen transport through POPC membrane using replica exchange transition interface sampling. *J. Phys. Chem. B.* 125:193–201.

30. Riccardi, E., O. Dahlen, and T. S. van Erp. 2017. Fast decorrelating Monte Carlo moves for efficient path sampling. *J. Phys. Chem. Lett.* 8:4456–4460.

31. Zhang, D. T., E. Riccardi, and T. S. van Erp. 2023. Path sampling with sub-trajectory moves. *J. Chem. Phys.* 158:024113.

32. Roet, S., D. T. Zhang, and T. S. van Erp. 2022. Exchanging replicas with unequal cost, infinitely and permanently. *J. Phys. Chem. A.* 126:8878–8886.

33. Davoudi, S., and A. Ghysels. 2023. Defining permeability of curved membranes in molecular dynamics simulations. *Biophys. J.* 122:1–10. https://doi.org/10.1016/j.bpj.2022.11.028.

34. Dellago, C., P. G. Bolhuis, and P. L. Geissler. 2002. Transition path sampling. *Adv. Chem. Phys.* 123:1.

35. Hummer, G. 2005. Position-dependent diffusion coefficients and free energies from Bayesian analysis of equilibrium and replica molecular dynamics simulations. *New J. Phys.* 7:34.

36. Ghysels, A., R. M. Venable, …, G. Hummer. 2017. Position-dependent diffusion tensors in anisotropic media from simulation: oxygen transport in and through membranes. *J. Chem. Theory Comput.* 13:2962–2976.

37. Krämer, A., A. Ghysels, …, R. W. Pastor. 2020. Membrane permeability of small molecules from unbiased molecular dynamics simulations. *J. Chem. Phys.* 153, 124107.

38. Dellago, C., P. G. Bolhuis, and D. Chandler. 1998. Efficient transition path sampling: application to Lennard-Jones cluster rearrangements. *J. Chem. Phys.* 108:9236–9245.

39. Hastings, W. 1970. Monte-Carlo sampling methhods using Markov chains and their applications. *Biometrika.* 57:97.

40. van Erp, T. S. 2007. Reaction rate calculation by parallel path swapping. *Phys. Rev. Lett.* 98, 268301.

41. Cabriolu, R., K. M. Skjelbred Refsnes, …, T. S. van Erp. 2017. Foundations and latest advances in replica exchange transition interface sampling. *J. Chem. Phys.* 147, 152722.

42. Dellago, C., P. G. Bolhuis, …, D. Chandler. 1998. Transition path sampling and the calculation of the rate constant. *J. Chem. Phys.* 108:1964–1977.

43. van Erp, T. S., M. Moqadam, …, A. Lervik. 2016. Analyzing complex reaction mechanisms using path sampling. *J. Chem. Theory Comput.* 12:5398–5410.

44. Plattner, N., J. D. Doll, …, J. E. Gubernatis. 2011. An infinite swapping approach to the rare-event sampling problem. *J. Chem. Phys.* 135, 134111.

45. Plattner, N., J. D. Doll, and M. Meuwly. 2013. Overcoming the rare event sampling problem in biological systems with infinite swapping. *J. Chem. Theory Comput.* 9:4215–4224.

46. Yu, T.-Q., J. Lu, …, E. Vanden-Eijnden. 2016. Multiscale implementation of infinite-swap replica exchange molecular dynamics. *Proc. Natl. Acad. Sci. USA.* 113:11744–11749.

47. Lu, J., and E. Vanden-Eijnden. 2019. Methodological and computational aspects of parallel tempering methods in the infinite swapping limit. *J. Stat. Phys.* 174:715–733.
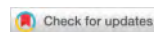
123

48. van Erp, T. S., and P. G. Bolhuis. 2005. Elaborating transition interface sampling methods. *J. Comput. Phys.* 205:157–181.

49. Moqadam, M., A. Lervik, …, T. S. van Erp. 2018. Local initiation conditions for water autoionization. *Proc. Natl. Acad. Sci. USA.* 115:E4569–E4576.

50. Lervik, A., E. Riccardi, and T. S. van Erp. 2017. PyRETIS: a well-done, medium-sized python library for rare events. *J. Comput. Chem.* 38:2439–2451.

51. Riccardi, E., A. Lervik, …, T. S. van Erp. 2020. PyRETIS 2: an improbability drive for rare events. *J. Comput. Chem.* 41:370–377.

52. Di, L., P. Artursson, …, K. Sugano. 2020. The critical role of passive permeability in designing successful drugs. *ChemMedChem.* 15: 1862–1874.

53. Choi, S.-H., S. Aid, and F. Bosetti. 2009. The distinct roles of cyclooxygenase-1 and -2 in neuroinflammation: implications for translational research. *Trends Pharmacol. Sci.* 30:174–181.

54. Novakova, I., E.-A. Subileau, …, W. Neuhaus. 2014. Transport rankings of non-steroidal antiinflammatory drugs across blood-brain barrier in vitro models. *PLoS One.* 9, e86806.

55. Jämbeck, J. P. M., and A. P. Lyubartsev. 2013. Exploring the free energy landscape of solutes embedded in lipid bilayers. *J. Phys. Chem. Lett.* 4:1781–1787.

56. Boggara, M. B., and R. Krishnamoorti. 2010. Partitioning of nonsteroidal antiinflammatory drugs in lipid membranes: a molecular dynamics simulation study. *Biophys. J.* 98:586–595.

57. Vanden-Eijnden, E., and M. Venturoli. 2009. Markovian milestoning with Voronoi tessellations. *J. Chem. Phys.* 130, 194101.

124

# 10

# PAPER II (PUBLISHED): PYRETIS 3: CONQUERING RARE AND SLOW EVENTS WITHOUT BOUNDARIES

D. T. Zhang and W. Vervust contributed equally to this work.
W. Vervust contributed to development of the PyRETIS 3 software package and writing of the manuscript.

**RESEARCH ARTICLE**

Journal of COMPUTATIONAL CHEMISTRY   WILEY

# PyRETIS 3: Conquering rare and slow events without boundaries

**Wouter Vervust**[1] | **Daniel T. Zhang**[2] | **An Ghysels**[1] | **Sander Roet**[3] |
**Titus S. van Erp**[2] | **Enrico Riccardi**[4]

[1]IBiTech–BioMMedA Group, Ghent University, Ghent, Belgium

[2]Department of Chemistry, Norwegian University of Science and Technology, Trondheim, Norway

[3]Department of Chemistry, Utrecht University, Utrecht, The Netherlands

[4]Department of Energy Resources, University of Stavanger, Stavanger, Norway

**Correspondence**
Enrico Riccardi, Department of Energy Resources, University of Stavanger, Stavanger, Norway.
Email: enrico.riccardi@uis.no

**Abstract**

We present and discuss the advancements made in PyRETIS 3, the third instalment of our Python library for an efficient and user-friendly rare event simulation, focused to execute molecular simulations with replica exchange transition interface sampling (RETIS) and its variations. Apart from a general rewiring of the internal code towards a more modular structure, several recently developed sampling strategies have been implemented. These include recently developed Monte Carlo moves to increase path decorrelation and convergence rate, and new ensemble definitions to handle the challenges of long-lived metastable states and transitions with unbounded reactant and product states. Additionally, the post-analysis software PyVisa is now embedded in the main code, allowing fast use of machine-learning algorithms for clustering and visualising collective variables in the simulation data.

**KEYWORDS**
kinetics, path sampling, PyRETIS, Python, rare event, slow event

## 1 | INTRODUCTION

The constant increase in high-performance computing (HPC) power enables molecular simulations to consider an increasing number of particles and a significantly longer simulated time. These hardware advancements have been complemented by the development of more efficient algorithms and software, substantially amplifying the effectiveness and predictive capacity of these methods. Despite these advancements, the study of rare transitions remains a computational challenge, as conventional simulations often capture insufficient transition events for statistical kinetic analysis.

For a numerical example, consider the dissociation rate of the drug molecule imatinib from the kinase protein ABL, which is approximately $10^{-3}$ s$^{-1}$.[1] Using a solvated simulation box of about 50,000 atoms, one can simulate up to 300 ns per day using a recent A100 GPU with an AMD EPYC CPU on the GROMACS MD software (version 2021.3).[2] Assuming dissociation events follow a Poisson process, one expects to observe an unbinding event after

9126 millennia of simulation time. Clearly, an increase of simulation speeds alone, albeit several orders of magnitude, will not suffice to extract kinetic information from biologically and chemically relevant systems.

In response, rare event simulations have emerged as a pivotal algorithmic advancement, enhancing the capabilities of molecular dynamics simulations for studying transition events.[3] In particular, the replica exchange transition interface sampling (RETIS) technique has proven to be one of the most accurate and efficient techniques for computing exact quantitative unbiased dynamical properties.[4] Yet, the necessity persists for intuitive and user-friendly software that can broaden the utilisation of these advanced simulation techniques, aiming to establish them as a standard tool accessible to non-specialists.

With this rationale, we developed PyRETIS, a Python library dedicated to efficiently simulating rare events in molecular systems based on the RETIS algorithm. Since its first release in 2017,[5] PyRETIS has undergone significant improvements, including extended features and algorithmic refinement in PyRETIS 2 in 2020.[6] To the best of our knowledge, PyRETIS is currently one of only two publicly available

Wouter Vervust and Daniel T. Zhang share the first authorship of the present work.

wileyonlinelibrary.com/journal/jcc   *J Comput Chem.* 2024;45:1224–1234.

126

path sampling codes capable of executing RETIS, the other being Open Path Sampling (OPS) code,[7,8] which was released in 2019.

PyRETIS, in its different releases, has been shown effective in multiple studies on a rather broad set of topics ranging from chemical reactions,[9–12] the adoption of a bacterial protein to DNA,[13] thin film breakage,[14,15] the solid-solid transition between the wurtzite and rock salt crystal structures,[16] and the oxygen permeation through membranes.[17] These comprehensive studies have served as a sturdy foundation, guiding the software's development over the past three years. A graphical representation of the latest works is included in Figure 1.

The latest iteration introduced in this article, PyRETIS 3, incorporates novel structural and algorithmic enhancements. Notably, this new release boasts an optimised architecture designed for parallel simulations and features an interface with machine learning algorithms, simplifying the post-analysis of simulation results. Its modular structure is visually depicted in the accompanying flowchart (Figure 2). A comprehensive discussion of the theory supporting the software's development can be found in the literature.[3,4,17,21–24]

## 2 | ALGORITHMIC DEVELOPMENTS

### 2.1 | PPTIS: Partial path transition interface sampling

In replica exchange transition interface sampling (RETIS), the sampling is conducted by generating a large number of trajectories that are accepted or rejected based on the Metropolis-Hastings law.[25,26] To use RETIS, a set of interfaces along a main collective variable has to be positioned, where each of these interfaces defines an ensemble. A trajectory belongs to the path ensembles $[i^+]$ when it starts at $\lambda_0$ ($= \lambda_A$, the boundary around stable state definition for the reactant state A), crosses the specific interface $\lambda_i$ and reaches either the product state at $\lambda_B$ or completely returns to the reactant state at $\lambda_0$, as visualized in Figure 3A. However, when the transition is not only rare, but also slow, as paths may get stuck in local metastable states, the average path lengths may extend beyond tens or hundreds of nanoseconds for some of the RETIS ensembles.



**FIGURE 1** Snapshots and descriptive representations of the latest works performed using the PyRETIS simulation library. (A) Thin film breakage[14,15] was investigated with force field dynamics using the GROMACS engine.[2] (B) The permeability of ibuprofen through a bilayer membrane[18] was estimated using the PyRETIS internal engine. (C) The increase in ergodicity by addition of the replica exchange move to PPTIS was demonstrated using a maze system.[18] (D) The electron transfer reaction between two ruthenium $(2+/3+)$ ions[15,19] was investigated with *ab initio* dynamics using the CP2K[20] engine.

**FIGURE 2**    Flowchart of PyRETIS 3 logic.



To address this limitation, we have implemented partial path transition interface sampling (PPTIS)[27] into the PyRETIS code. With this method, the path ensembles $[i^{\pm}]$ contain paths that cross $\lambda_i$ and start and end on neighbouring interfaces $\lambda_{i-1}$ or $\lambda_{i+1}$, as visualized in Figure 3B. From the paths, the local crossing probabilities $p_i^{\pm}$ and $p_i^{\mp}$ can be determined, and with a recursive relation,[27] the global crossing probability $P_A(\lambda_B|\lambda_A)$ can be estimated based on the assumption of memory loss.

In PyRETIS 3, users can opt for a PPTIS simulation with a simple keyword in the input.rst file, as described in the online documentation. In the summary file pathensemble.txt, the path labels of the $[i^{\pm}]$ ensemble can now take on four different path-types for accepted paths: LML or LMR (as in a TIS simulation), and also RMR or RML (Figure 3B). Here, 'L', 'M', and 'R' denote the left, middle, and right interfaces within a designated path ensemble. In post-analysis, the ratio of different path-types determines local crossing probabilities,

**FIGURE 3** (A) Trajectories of a TIS path ensemble $[i^+]$ start from the reactant interface $\lambda_A$ (L), cross the specific interface $\lambda_i$ (M) and reach either the product interface $\lambda_B$ (R) or completely return to $\lambda_A$. (B) Trajectories of a PPTIS path ensemble $[i^\pm]$ start from either the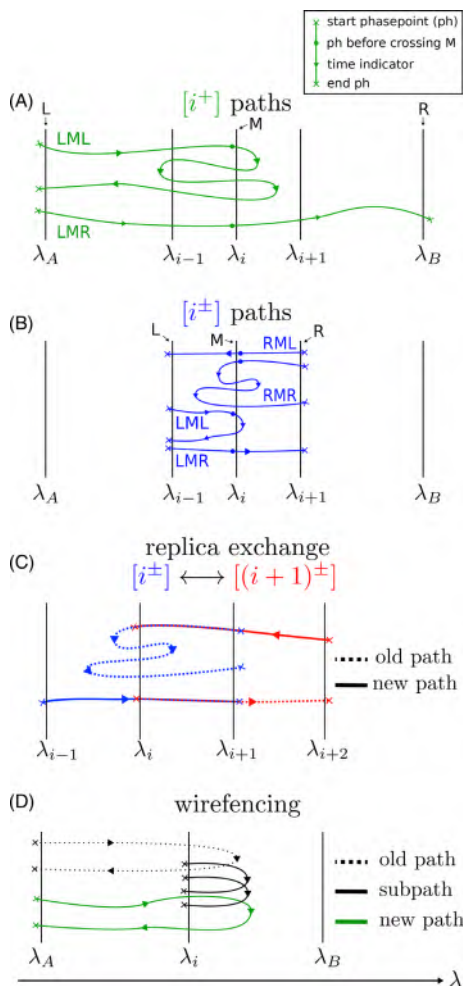 left-neighbouring (L) or right-neighbouring (R) interface of $\lambda_i$ (M), after which they cross M, and return to either L or R. (C) Schematic representation of a replica exchange move in a REPPTIS simulation. The RMR path of $[i^\pm]$ (blue dotted) is extended backward in time to create a new RML path of $[(i+1)^\pm]$ (solid red). Likewise, the LMR path of $[(i+1)^\pm]$ (red dotted) is extended backwards in time to create a new LMR path of $[i^\pm]$ (blue solid). (D) Schematic representation of a WF move in a RETIS simulation. A new path (solid green) is generated from an old path (dotted black) via a chain of subtrajectories (solid black), which can significantly improve path decorrelation compared to the standard shooting move.

for example, $p_i^\pm = \text{LMR}/(\text{LMR} + \text{LML})$ denotes the local probability of reaching the right interface $\lambda_{i+1}$ rather than the left interface $\lambda_{i-1}$, given that a path has crossed the middle interface $\lambda_i$ after having crossed the left interface $\lambda_{i-1}$.

## 2.2 | REPPTIS: Replica exchange partial path transition interface sampling

Akin to RETIS enhancing the transition interface sampling (TIS)[28] efficiency through replica exchange,[29–31] similar advancements in PPTIS can lead to the REPPTIS method.[18] The inclusion of replica exchange in REPPTIS relies on MC moves designed to exchange path segments among adjacent path ensembles. These segments are then extended until they meet the acceptance criteria of a neighbouring ensemble, an example of which is shown in Figure 3C. As shown in Reference 18, these swaps notably facilitate the exploration in regions partitioned by energetic or kinetic barriers orthogonal to the order parameter $\lambda$, thereby substantially improving the sampling's ergodicity

The acceptance of a swap hinges on the pairing of path-types (LMR, LML, RML, or RMR) and the selected propagation directions (forward or backward in time). For example, an LMR path in the $[2^\pm]$ ensemble can be extended into $[3^\pm]$ by forward propagation in time, while backward propagation would push the path towards the $[1^\pm]$ ensemble. If propagation directions are incompatible with the path-types, the replica exchange move is immediately rejected without requiring MD steps, resulting in an SWD entry ('Swap Wrong Direction') in the pathensemble.txt files of the relevant ensembles. Next, if propagation directions are compatible with path-types, path segments are extended until they cross either the left or right interface of the new ensemble.

Contrary to swaps between the RETIS $[i^+]$ ensembles, swaps between REPPTIS ensembles require MD integration in the extension step. In RETIS, only the swap between $[0^-]$ and $[0^+]$ require MD integration. Fortunately, in both RETIS and REPPTIS, any required MD integration commences only upon confirmation of swap acceptance. Consequently, the CPU time spent on MD integration is almost never wasted. The only exception is when paths exceed the user-specified maximum path length. Such occurrences should be rare as (RE)PPTIS is designed to reduce path lengths. Frequent occurrences suggest that the user should either add more interfaces, increase the maximum path length parameter, or change the order parameter definition. An example PyRETIS 3 input file running REPPTIS with a certain swap frequency is shown in Figure 4. The user can refer to the online documentation for detailed instructions on setting up a REPPTIS simulation.

## 2.3 | Improved stone skipping & web throwing

To increase the simulation efficiency, advanced shooting moves in the form of stone skipping (SS) and web throwing (WT)[23] were already implemented in PyRETIS 2. Compared to standard shooting

```
Simulation
----------
task = repptis
steps = 10000
interfaces = [-0.5, -0.3, 0.0, 0.3, 0.5]
permeability = True
zero_left = -0.75


TIS settings
------------
freq = 0.0
maxlength = 100000


RETIS settings
--------------
swapfreq = 0.2
swapsimul = True
```

**FIGURE 4**   New keywords introduced to perform a REPPTIS simulation. Replica exchange moves are introduced by setting the 'task' keyword to repptis. Typically, one will not perform time-reversal moves in the PPTIS framework, which is accomplished by setting the 'freq' keyword of TIS settings to 0.0. The relative frequency of shooting moves to replica exchange moves is set by the 'swapfreq' keyword of the RETIS settings. In this example, 20 % of the moves will perform replica exchange. This input excerpt is tailored for a permeability simulation, as the 'permeability' keyword is set to True. The positioning of the accompanying $\lambda_{-1}$ interface is done by setting the 'zero_left' keyword to $-0.75$.

methods,[32] both SS and WT are more cost-efficient per MD step. They achieve this by leveraging a sequence of intermediate short paths (subpaths) that minimise correlations between previously existing and newly generated paths. After a number of subpaths have been completed, the last one is extended to become a new full path. Although SS and WT yield an increase in efficiency of more than one order of magnitude in case studies,[23] their practical implementation can be impeded when linked with external MD engines or when the calculation of the order parameter is expensive.

More specifically, new SS/WT subpaths must be launched from a randomly perturbed phase point (time slice) of the previous path/subpath that establishes a crossing with some relevant interface (like the main ensemble's interface that always needs to be crossed). Typically, only the velocities are changed in this perturbation, but after this process, it should be verified that the perturbed phase point is still a crossing point: the other side of the interface should be reached after a time step forward or backward in time. If this is not the case, the move is not directly rejected, but new velocities should be generated. Especially if the time step considered by the RETIS algorithm actually consists of several MD steps, this requirement is not always easy to fulfil and multiple velocity randomization attempts can be required.

If a RETIS time step is a single MD step (the subcycle is set to 1), Reference 15 proposed two strategies to minimise the computational cost for doing the one-step crossing test. First, the velocity-Verlet integrator[15,33] can be reformulated to predict the next configuration point without the

associated expensive force calculation. Therefore, given a configuration-based order parameter and a relatively cheap velocity generation, the cost of testing reduces considerably such that the one-step crossing requirement can quickly be achieved. The second strategy is to modify the velocity generation procedure to increase the likelihood of pushing the next step across the interface. For example, for an $N$ particle system, certain velocities can be kept or reversed instead of letting all $3N$ velocity components be regenerated from a Maxwell-Boltzmann distribution.

However, in the case that one MD engine call implies performing several MD steps (i.e., subcycle is set to 10-2000), and/or when the order parameter calculation is expensive, the increased testing cost for the one-step crossing condition reduces the SS and WT computational efficiency. We therefore developed a third type of advanced shooting move that does not require the one-step crossing condition, which is implemented in PyRETIS 3 and discussed in the next section.

## 2.4 | Wire fencing

To circumvent the issues related to the one-step crossing condition, we formulated a third advanced shooting move within the subtrajectory family that we named "wire fencing" (WF).[15] In comparison to SS and WT, a new WF subpath can be launched from a configuration point of the previous path or subpath that lies between the path ensemble's interface $\lambda_i$ and $\lambda_{cap}$, a user-defined *cap interface*: $\lambda_i < \lambda_{cap} \le \lambda_B$. It is therefore not restricted to crossing points of specific interfaces. In the cases of potential energy surfaces that have gradual regions close to state $B$, the shooting point selection can be restricted to only occur in steep regions by a suitable $\lambda_{cap}$ placement. A schematic representation of the WF move is given in Figure 3D.

The WF move, like all advanced shooting moves, are best combined with the *high-acceptance* technique.[15,23] The high-acceptance scheme alters the sampled path distribution, a change that can be precisely adjusted through a reweighting scheme in the post-simulation analysis. In this specific formulation, the Metropolis-Hastings acceptance criteria solely reject paths that both commence and culminate at the interface $\lambda_B$.

To use SS, WT or WF in PyRETIS 3, the user can select, for each ensemble, the type of shooting-moves, the number of subpaths, the positions of the additional interfaces (like the cap-interface in WF), and whether to adopt the high-acceptance sampling scheme.

## 2.5 | The $\lambda_{-1}$ interface

In both RETIS and REPPTIS simulations, the flux term is computed from the average path lengths of the $[0^-]$ and $[0^+]$ ensembles. The former ensemble, $[0^-]$, is conventionally defined by a single interface and obeys different sampling rules compared to the other ensembles. In some modelling situations, it could be advantageous to restrict the $[0^-]$ ensemble within two interfaces. A first example is when the reactant state $A$ is unbound to the left, such as when it signifies an infinite reservoir or when it is barrierlessly linked to state $B$ through periodic boundary conditions.

In such cases, while the rate might not be well defined, other dynamical properties—such as the permeability in a membrane system—are still ascertainable. A second example is when the reactant state $[0^-]$ extends over a finite, but very large collective variable $\lambda$ range. Here, it may be more strategic to focus on sampling the region closer to $\lambda_0$. In Reference 34, an additional interface $\lambda_{-1}$ was introduced to the left of $\lambda_0$, and the new ensemble $[0^{-\prime}]$ was defined as the path ensemble containing trajectories that start and stop on $\lambda_0$ or $\lambda_{-1}$. Consequently, the region $\lambda < \lambda_{-1}$ is never sampled, hence avoiding periodic boundary crossings or waste of computer time in a non-relevant region of phase space.

A theoretical complication is that an adaptation is needed to match the $[0^-]$ and $[0^+]$ ensemble. Indeed, the number of paths that may be cut out from a very long equilibrium MD simulation and that belong to the standard $[0^-]$ is equal to those of $[0^+]$, that is, $N_{[0^-]} = N_{[0^+]}$. The $\lambda_{-1}$ interface increases the number of paths, however, as the trajectory is cut more often, and $N_{[0^{-\prime}]} > N_{[0^+]}$. To get the correct permeability, a correction $\xi$ is therefore needed,[34] where $\xi = N_{[0^+]}/N_{[0^{-\prime}]}$. Fortunately, the factor $\xi$ can be readily evaluated in the post-processing by PyRETIS-3.

## 2.6 | Mirror move and target-swap move

The permeability coefficient can directly be computed in PyRETIS 3. It can be obtained by considering the rate of transition while following a single permeant, referred to as the target, from left to right through a given region (e.g., a membrane). The presence of other permeants in the system affects the rate, as they are part of the surrounding environment. Transitions from right to left are prevented due to the presence of the $\lambda_{-1}$ interface. However, for the sake of sampling efficiency, it would be beneficial to incorporate the statistics of other transitioning permeants and utilise transitions in both directions, including from right to left, whenever the membrane is (statistically) symmetric. This potential for sampling efficiency gain is achieved through the target-swap move and the mirror move, respectively.

The mirror move involves mirroring the $z$ coordinates of all particles of all time slices of the previous path across an $xy$-mirror plane located between two periodic images of the membrane. Additionally, the $z$-component of the particle velocities is flipped. However, it is important to note that the mirror move must be accompanied by setting the $\lambda_{-1}$ interface at an equal distance away from the mirror plane in the water slab as the $\lambda_0$ interface, but on the opposite side. The mirror move solely applies to the $[0^{-\prime}]$ ensemble and has the consequence that a previous path ending at $\lambda_{-1}$, will end at $\lambda_0$ after this move. It is worth noting the distinction from the time-reversal move, as the mirror move would lead to a switch to the opposite interface even if the previous path started and ended at the same side. As a result of the new path now ending at $\lambda_0$, it can be successfully swapped with a $[0^+]$ path in the next MC step.

For code-technical reasons, PyRETIS-3 implements the mirror move in a slightly different but equivalent manner. Instead of changing the particle coordinates and velocities directly, PyRETIS-3 modifies the definition of the reaction coordinate (RC) by using the $z$ of the target's position, and mirroring it across a mirror plane. This alternative implementation facilitates the integration of PyRETIS-3 with external molecular dynamics (MD) engines, which may have different methods for altering coordinates and velocities. The flag indicating the sign of $z$ to be used in the RC definition is also exchanged during the replica exchange moves of the RETIS algorithm.

The target-swap move also plays a crucial role in improving sampling efficiency. This move randomly selects a permeant to be considered as the new target from all the time slices of the previous path. It considers permeants within the $[0^{-\prime}]$ boundaries, namely $\lambda_{-1}$ and $\lambda_0$, at that time slice. Once a random time slice and permeant pair is chosen, the trajectory is traced both forward and backward in time until the new target reaches the boundaries. This process may involve extending or truncating the old trajectory along each time direction. Ultimately, the final trajectory is accepted or rejected based on a Metropolis decision, ensuring that detailed balance is maintained.

It is worth noting that even if the mirror move and the target swap move operate exclusively in the minus ensemble ($[0^{-\prime}]$), the enhanced exploration trickles down to the other ensembles through replica exchange swaps, improving the sampling across all ensembles. The effectiveness of both moves was distinctly demonstrated in a two-channel system where it significantly enhanced sampling in the collective variable space orthogonal to the reaction coordinate.[34]

## 3 | POST-PROCESSING

### 3.1 | Error analysis

Calculating statistical errors within a RETIS simulation is a non-trivial task due to the various types of correlation. In contrast to TIS, where path simulations are independent and standard error propagation rules can be used, RETIS introduces additional complexities. Although block averages can address correlations within a path ensemble, correlations extend across different path ensembles in RETIS due to path swapping, which results in shared data between the ensembles. The involvement of the $[0^+]$ path ensemble in both the flux and crossing probability calculations, along with the expected correlation between path length and reaction progress, further hinders assuming error independence. Additionally, computing the total crossing probability using the weighted histogram analysis method (WHAM)[21,35,36] also increases the difficulty for dealing with correlations as it utilises overlaps in path ensembles to improve the accuracy.

The implementation of standard block averaging poses practical challenges. This issue is evident in REPPTIS, where $p_i^{\pm}$ and $p_i^{\mp}$ may be based on different numbers of sampled trajectories, rendering an absolute block length unsuitable. One potential solution is to adapt the block length to the relative size of the different data sets. For example, in a RETIS/REPPTIS analysis, the first 10% of each data file could be utilised to compute a rate. Subsequently, data between 10% and 20% from all files are employed to calculate a new rate, continuing this process until ten nearly independent estimates are obtained.

These rates can then be used to compute the standard deviation and, ultimately, the error.

In a practical implementation, the described approach can still be viable if the number of blocks (10 in the aforementioned example) is predetermined and fixed. However, block averaging techniques often demonstrate significant fluctuations in computed errors due to the arbitrary choice of block length or, equivalently, the number of blocks. It is therefore recommended to conduct error analysis using a range of block lengths. By evaluating the computed errors across different block lengths and observing the graph of computed error versus block length $m$, one can identify a plateau region where the errors stabilise. This specific region is usually identified for sufficiently large $m$ values, where the blocks can be deemed uncorrelated yet remain small enough to maintain a substantial number of blocks. Taking the average of the computed errors within this plateau region is considered good practice, as it provides a more reliable estimate of the statistical uncertainty. Yet, partitioning all data files into many sets of relative blocks is cumbersome and time consuming.

We have, therefore, taken a more practical yet sound approach, that we refer to as recursive block errors, which surprisingly has not been widely mentioned as an alternative to standard block averaging. The recursive block errors approach is based on a single data file containing the running estimate as a function of performed MC cycles of the property under investigation, such as the rate, flux, crossing probability, or specific path ensemble properties like local crossing probabilities and path lengths. These running estimates are standard outputs from PyRETIS simulations.

Suppose the RETIS simulation consists of $N$ MC cycles, where a cycle typically implies the update of each path ensemble by a MC move. Let $k[n]$ with $1 \leq n \leq N$ be the running estimate of the rate after $n$ cycles have been completed. Given a block length of $m$, there are $M = \text{int}(N/m)$ blocks. We now introduce the *recursive-block* value $k_j^r$ for each block $j = 1,2,...,M$. These values are implicitly defined by ensuring that the mean of the initial $j$ recursive-block values matches the running estimate after $jm$ MC cycles:

$$k[jm] = (k_1^r + k_2^r + k_3^r + \cdots + k_j^r)/j \tag{1}$$

which leads to a recursive relation for $k_j^r$:

$$k_j^r = \begin{cases} jk[jm] - (k_1^r + k_2^r + k_3^r + \cdots + k_{j-1}^r) & \text{if } j > 1 \\ k[m] & \text{if } j = 1 \end{cases} \tag{2}$$

From this, given a specific block length $m$, the following non-recursive relation can be extracted

$$k_j^r = jk[jm] - (j-1)k[(j-1)m] \tag{3}$$

These block values $k_j^r$ are then used to compute the standard deviation between them and ultimately the estimated error for this particular block size $m$, just like is done with standard block averaging.

The great advantage of this simple relation is that the post-analysis can start using just the running estimate $k[n]$ with $1 \leq n \leq N$,

which subsequently can be reused to compute $k_j^r$ for a large set of block sizes $m$. Consequently, this method enables an efficient computation of the estimated error concerning block size, facilitating the extraction of a final robust error estimate by averaging specifically within the identified plateau region. In the appendix, we show that this approach converges to the same error estimates as when standard-block averages are used.

## 3.2 | Permeability

It is challenging to compute the permeability $P$ of permeants through a barrier (e.g., through a membrane), when a high free energy barrier needs to be overcome or when the permeants are trapped in a free energy well for an extended time.[37] An example of the former is the water permeation through phospholipid membranes,[38–40] and an example of the latter is the transport of oxygen molecules through membranes.[41,42] In Reference 34, it was shown that the permeability through a region can be estimated from the RETIS crossing probability $P_A(\lambda_B|\lambda_A)$. As an additional quantity for $P$, the average time that a path spends in a reference interval in the $[0^-]$ or $[0^{-\prime}]$ ensemble needs to be calculated. Moreover, the $\xi$ factor (see subsection on the $\lambda_{-1}$ interface) also needs to be evaluated when a $\lambda_{-1}$ interface is used. These two extra quantities were implemented in a straightforward way in the PyRETIS post-processing analysis tools.

With the PyRETIS implementation in Reference 34, a series of 1D examples of barriers with varying height and viscosity settings were tested. In addition, the permeation of particles through a 2D membrane with two distinct permeation pathways was successfully mapped out, where the use of the target-swap move and mirror move enhanced the sampling efficiency. When the permeation event is not only rare but also slow, the permeability can be assessed with the PPTIS or REPPTIS methodology.[18] The (RE)PPTIS implementation was used to investigate the permeation through a 2D maze, representing a membrane with two permeation pathways, where one has an entropic barrier and the other an energetic barrier. The computed (RE)PPTIS permeabilities were benchmarked against the RETIS permeabilities, where the replica exchange moves in REPPTIS could significantly reduce the effect of memory-loss in the partial path sampling method.

## 3.3 | PyVisA

PyVisA,[43] a post-processing visualisation and analysis tool, has been integrated into the PyRETIS 3 release, with its own executable and optional graphical user interfaces. The library permits to directly load and visualise simulation outcomes, reducing the data handling time and costs, as data can be checked and compressed remotely and visualised locally. The simulation results can be displayed as a whole or in sections. Different collective variables, simulation time, trajectory status, trajectory number can be directly grouped for different ensembles supporting the analysis and visualisation of simulation results. A customizable interface to improve the appearance of the reports has

been constructed and an interactive interface implemented. Furthermore, the selected data subset can be saved in ".hdf5", simple ".txt", or with a ".json" format. Images can also be exported directly as ".png". To provide further visualisation customizability, PyVisA can also output a minimalistic python script to regenerate the selected image from selected data. Clearly, the script can be then re-adapted to the best user convenience. If trajectory data ($x$-, $y$-, $z$-positions) are stored, each data-point of the displayed trajectory is linked by PyVisA to the original source, and a molecular 3D visualisation can be launched showing the molecular structure at the points of interest. The selected data can then be fed directly to a set of machine learning approaches such as clustering algorithms,[44] random forests,[45] calculation of the Pearson correlation matrix of coefficients,[46,47] and so forth. These approaches are provided by the scikit-learn python package[48] and the sampling data is internally wrangled such that it can directly be fed to these algorithms. PyVisA has been constructed to be executed independently while simulations are ongoing, allowing for intermediate descriptions and visualisations of the simulation outcomes. The new

panels to access the interactive functionality and the integrated machine learning modules are shown in Figure 5.

## 4 | OTHER UPDATES

### 4.1 | Code structure

In PyRETIS 2 users could define different settings for different ensembles, to allow for better customization in the design of sampling problems. In PyRETIS 3, the internal representation of ensembles has also been subdivided. Each ensemble is represented by an independent object, allowing multiple simulations to be performed contemporaneously. Each ensemble can also be, in principle, executed with a different external engine. The trajectory management has been optimised, where trajectory cleaning is now performed after every ensemble move rather than waiting for an entire MC cycle to finish. This significantly reduces storage space for large systems using many ensembles.
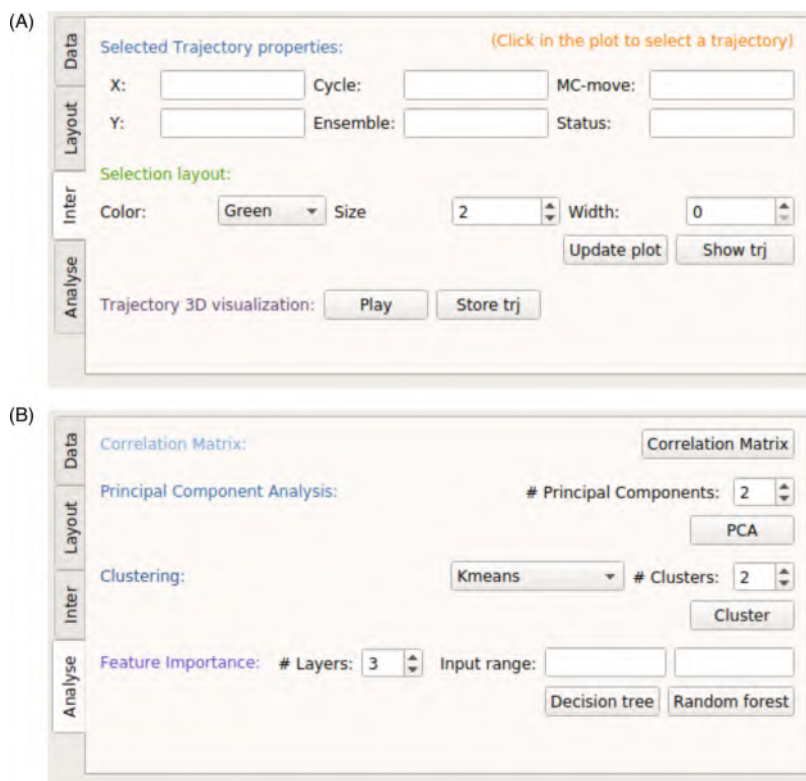


**FIGURE 5** PyVisA interactive visualisation panel (A) and PyVisA machine learning model control panel (B).

The new code structure allows for effective parallel simulations to be performed in different ensembles, potentially speeding up sampling performance by optimising cluster usage. An effective simulation handler procedure to advantageously execute multiple parallel simulations has been recently published by Roet et al.[49] and it can be expected to be included in forthcoming software releases.

## 4.2 | Interface with the external engine Amsterdam modeling suite

The PyRETIS code has been developed with the intent to use external engines (e.g., GROMACS,[2] LAMMPS,[50] CP2K,[20] OpenMM[51]) and to facilitate the implementation of new ones. In collaboration with SCM (Software for Chemistry & Materials), we have established an interface between PyRETIS and their Amsterdam Modeling Suite (AMS).[52] This development allows PyRETIS 3 to utilise AMS as its MD engine, significantly expanding the range of dynamic simulations it can perform. The extended capabilities include the ability to conduct *ab initio* MD based on the Amsterdam Density Functional (ADF) package, ReaxFF,[53] and Density Functional Tight Binding (DFTB).[54]

These functionalities are already implemented and are currently available in the development branch of the PyRETIS code. At the time of submitting this paper, the documentation and unit testing for this feature were not fully completed. Nevertheless, the functionalities are fully operational. We encourage users interested in utilising this functionality to reach out to the PyRETIS developers for support and guidance on its usage.

## 5 | USER GUIDELINES

To maximize the efficiency of the sampling strategies offered by PyRETIS, we provide some guidelines on which simulation technique is best suited for which applications. In general, RETIS is a very efficient sampling strategy for rare events, which are characterized by long waiting periods followed by the actual transition that happens very quickly. REPPTIS, on the other hand, is designed to tackle slow events. A slow event can still be characterized by a long waiting period, after which the transition happens in a sluggish fashion, where a transition spends considerable time in metastable states along the reactive pathway.

In more technical terms, RETIS is most suitable for reaction mechanisms that are well-described by overcoming a single, large energetic and/or entropic barrier. If the reaction mechanism includes metastable states, a collection of RETIS simulations can still be used if the location of the free energy wells are well-known beforehand. A Markov state model can then be built from the collection of RETIS simulations, from which the global transition rate can be estimated. Examples of this include the permeation of small molecules (methanol, ethanol, etc.) through nanoporous materials or phospholipid bilayers. When the system becomes more complex, the optimal reaction coordinate quickly becomes elusive due to the increasing number of transition

states and metastable states. Consequently, (unknown) metastable states emerge along (unknown) orthogonal degrees of freedom concerning the selected reaction coordinate, for which REPPTIS is better suited. Examples include the association and dissociation of drug molecules to proteins, the permeation of small drug molecules through phospholipid bilayers, protein-protein interactions, and so forth.

## 6 | USER SUPPORT

The code has been updated with the latest python libraries and its dependencies on external packages have been minimized in order to increase its maintainability. Our software development has been inspired by the FAIR software principles. The code is fully open source and shared on GitLab and GitHub.[55] Documentation, tests, and examples are constantly updated to simplify the code usability by external users.

## 7 | CONCLUSIONS AND FUTURE WORK

We have released the third version of the PyRETIS code, which includes algorithmic developments that can greatly improve sampling efficiency. In particular, partial path transition interface sampling (PPTIS), replica exchange partial path transition interface sampling (REPPTIS), and advanced shooting moves (wire fencing, mirror move, and target-swap move) were implemented and previous implementations of the stone skipping and web throwing moves were improved. Additionally, PyRETIS 3 solidifies the $\lambda_{-1}$ interface feature based on a new well defined and relevant path ensemble $[0^{-\prime}]$, which improves the purely pragmatic implementation of the $\lambda_{-1}$ interface option in PyRETIS 2. Within the $[0^{-\prime}]$ ensemble, two additional terms are calculated: the average time that a path in this ensemble spends in a reference region and the $\xi$ factor. These two values, together with the computed RETIS crossing probability, allow one to compute the permeability coefficient in bounded and unbounded systems. The approach does also not require any additional flux calculations, which was the pragmatic solution suggested in the PyRETIS 2 for bounded systems. A direct approach to compute the permeability in unbounded systems was not available in PyRETIS 2.

In the post-processing phase of the code, we enhanced our error analysis by employing *recursive block* analysis. This method is particularly effective in handling intricate correlations, especially in scenarios where the final computed properties stem from a series of separate yet interdependent simulations, a characteristic notably pertinent to RETIS and REPPTIS. We have then further improved and formally included PyVisA as a part of the PyRETIS 3 release. The library provides an analysis and visualisation toolkit to facilitate the study of simulation output. In particular, the library prepares the data such that machine learning approaches (e.g., random forest, clustering) can be directly applied, providing enhanced insight on the results.

Further software development will focus on code parallelization. The recent introduction of ∞RETIS,[49,56] a novel RETIS variant

integrating asynchronous replica exchange with infinite swaps, represents a crucial advancement in line with this objective. This innovation effectively resolves the challenge posed by the imbalanced CPU costs associated with the (advanced) shooting moves, stemming from varying path lengths.

In the domain of force field development and data analysis, the synergy between path sampling and machine learning will continue to advance. For instance, leveraging RETIS simulations at the Ab Initio MD level can yield a pertinent training set of configuration points crucial for establishing a reactive force field through neural networks[57–59] tailored for specific chemical reactions. Once the force field is established, it can be reintegrated into the RETIS simulation, offering unprecedented convergence and reliability. Additionally, mechanistic analysis through committor analysis[60,61] and the predictive capacity analysis[10,21] is anticipated to provide deeper insights with the aid of machine learning tools. The objective is not solely to comprehend the occurrence of rare events like chemical reactions but also to understand methods for enhancing, directing, or impeding them. The PyRETIS development team aims to propel these advancements by offering new open-source computer codes that are publicly accessible, intuitive, and instrumental in tackling complex applications while advancing quantitative path sampling algorithms.

## ACKNOWLEDGMENTS

## DATA AVAILABILITY STATEMENT

The data that support the findings of this study are openly available in pyretis at https://gitlab.com/pyretis/pyretis.git.

## ORCID

*Wouter Vervust* https://orcid.org/0000-0002-8714-3017
*Daniel T. Zhang* https://orcid.org/0000-0002-0296-0860
*An Ghysels* https://orcid.org/0000-0003-0015-2605
*Sander Roet* https://orcid.org/0000-0003-0732-545X
*Titus S. van Erp* https://orcid.org/0000-0001-6600-6657
*Enrico Riccardi* https://orcid.org/0000-0003-1890-7113

## REFERENCES

[1] A. Lyczek, B.-T. Berger, A. M. Rangwala, Y. Paung, J. Tom, H. Philipose, J. Guo, S. K. Albanese, M. B. Robers, S. Knapp, et al., *Proc Natl Acad Sci* **2021**, *118*, e2111451118.

[2] M. J. Abraham, T. Murtola, R. Schulz, S. Páll, J. C. Smith, B. Hess, E. Lindahl, *SoftwareX* **2015**, *1–2*, 19.

[3] T. S. van Erp, *Solvent effects on chemistry with alcohols*. Ph.D. thesis, Universiteit van Amsterdam, Netherlands **2003**.

[4] T. S. van Erp, *Phys. Rev. Lett.* **2007**, *98*, 268301.

[5] A. Lervik, E. Riccardi, T. S. van Erp, *J. Comput. Chem.* **2017**, *38*, 2439.

[6] E. Riccardi, A. Lervik, S. Roet, O. Aarøen, T. S. van Erp, *J. Comput. Chem.* **2020**, *41*, 370.

[7] D. W. H. Swenson, J.-H. Prinz, F. Noe, J. D. Chodera, P. G. Bolhuis, *J. Chem. Theory Comput.* **2019**, *15*, 813.

[8] D. W. H. Swenson, J.-H. Prinz, F. Noe, J. D. Chodera, P. G. Bolhuis, *J. Chem. Theory Comput.* **2019**, *15*, 837.

[9] M. Moqadam, E. Riccardi, T. T. Trinh, A. Lervik, T. S. van Erp, *Phys. Chem. Chem. Phys.* **2017**, *19*, 13361.

[10] M. Moqadam, A. Lervik, E. Riccardi, V. Venkatraman, B. K. Alsberg, T. S. van Erp, *Proc. Natl. Acad. Sci. U. S. A.* **2018**, *115*, E4569.

[11] S. Roet, C. D. Daub, E. Riccardi, *J. Chem. Theory Comput.* **2021**, *17*, 6193.

[12] C. D. Daub, E. Riccardi, V. Hänninen, L. Halonen, *PeerJ Phys Chem* **2020**, *2*, e7.

[13] E. Riccardi, E. C. van Mastbergen, W. W. Navarre, J. Vreede, *PLoS Comput. Biol.* **2019**, *15*, e1006845.

[14] O. Aarøen, E. Riccardi, T. S. van Erp, M. Sletmoen, *Colloids Surf., A* **2022**, *632*, 127808.

[15] D. T. Zhang, E. Riccardi, T. S. van Erp, *J Chem Phys* **2023**, *158*, 024113.

[16] A. Lervik, I.-H. Svenum, Z. Wang, R. Cabriolu, E. Riccardi, S. Andersson, T. S. van Erp, *Phys. Chem. Chem. Phys.* **2022**, *24*, 8378.

[17] E. Riccardi, A. Krämer, T. S. van Erp, A. Ghysels, *J Phys Chem B* **2020**, *125*, 193.

[18] W. Vervust, D. T. Zhang, T. S. van Erp, A. Ghysels, *Biophys. J.* **2023**, *122*, 2960.

[19] A. Tiwari, B. Ensing, *Faraday Discuss.* **2016**, *195*, 291.

[20] J. Hutter, M. Iannuzzi, F. Schiffmann, J. VandeVondele, *WileyWIREs Comput Mol Sci* **2014**, *4*, 15.

[21] T. S. van Erp, M. Moqadam, E. Riccardi, A. Lervik, *J. Chem. Theory Comput.* **2016**, *12*, 5398.

[22] T. S. van Erp, T. P. Caremans, C. E. A. Kirschhock, J. A. Martens, *Phys. Chem. Chem. Phys.* **2007**, *9*, 1044.

[23] E. Riccardi, O. Dahlen, T. S. van Erp, *J. Phys. Chem. Lett.* **2017**, *8*, 4456.

[24] T. S. van Erp, *Epl* **2023**, *143*, 30001.

[25] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, E. Teller, *J. Chem. Phys.* **1953**, *21*, 1087.

[26] W. K. Hastings, *Biometrika* **1970**, *57*, 97.

[27] D. Moroni, P. G. Bolhuis, T. S. van Erp, *J. Chem. Phys.* **2004**, *120*, 4055.

[28] T. S. van Erp, D. Moroni, P. G. Bolhuis, *J. Chem. Phys.* **2003**, *118*, 7762.

[29] R. H. Swendsen, J. S. Wang, *Phys. Rev. Lett.* **1986**, *57*, 2607.

[30] E. Marinari, G. Parisi, *Europhys Lett.* **1992**, *19*, 451.

[31] Y. Sugita, Y. Okamoto, *Chem. Phys. Lett.* **1999**, *314*, 141.

[32] C. Dellago, P. G. Bolhuis, D. Chandler, *J. Chem. Phys.* **1998**, *108*, 9236.

[33] D. Frenkel, B. Smit, *Understanding Mol. Simul*, 2nd ed., Academic Press, San Diego, CA **2002**.

[34] A. Ghysels, S. Roet, S. Davoudi, T. S. van Erp, *Phys. Rev. Res.* **2021**, *3*, 033068.

[35] A. M. Ferrenberg, R. H. Swendsen, *Phys. Rev. Lett.* **1989**, *63*, 1195.

[36] J. Rogal, W. Lechner, J. Juraszek, B. Ensing, P. G. Bolhuis, *J. Chem. Phys.* **2010**, *133*, 174109.

[37] S. Davoudi, A. Ghysels, *J. Chem. Phys.* **2021**, *154*, 054106.

[38] A. Krämer, A. Ghysels, E. Wang, R. M. Venable, J. B. Klauda, B. R. Brooks, R. W. Pastor, *J. Chem. Phys.* **2020**, *153*, 124107.

[39] R. M. Venable, A. Krämer, R. W. Pastor, *Chem. Rev.* **2019**, *119*, 5954.

[40] A. Ghysels, A. Krämer, R. Venable, W. Teague, E. Lyman, K. Gawrisch, R. W. Pastor, *Nat. Commun.* **2019**, *10*, 5616.

[41] A. Ghysels, R. M. Venable, R. W. Pastor, G. Hummer, *J. Chem. Theory Comput.* **2017**, *13*, 2962.

[42] O. De Vos, R. M. Venable, T. Van Hecke, G. Hummer, R. W. Pastor, A. Ghysels, *J. Chem. Theory Comput.* **2018**, *14*, 3811.

[43] O. Aarøen, H. Kiær, E. Riccardi, *J. Comput. Chem.* **2021**, *42*, 435.

[44] R. Xu, D. Wunsch, *IEEE Trans. Neural Netw* **2005**, *16*, 645.

[45] T. K. Ho, Random decision forests, Proceedings of 3rd International Conference on Document Analysis and Recognition, Montreal, QC, Canada, pp. 278–282 vol.1.

[46] K. Pearson, *Proc R Soc London* **1895**, *58*, 240.

[47] O. Dahlen, A. Lervik, O. Aaroen, R. Cabriolu, R. Lyng, T. S. van Erp, *Comput. Appl. Eng. Educ.* **2020**, *28*, 779.

[48] E. Jones, T. Oliphant, P. Peterson, et al., SciPy: Open source scientific tools for Python. http://www.scipy.org **2001**.

[49] S. Roet, D. T. Zhang, T. S. van Erp, *J Phys Chem A* **2022**, *126*, 8878.

[50] S. Plimpton, *J. Comput. Phys.* **1995**, *117*, 1.

[51] P. Eastman, M. S. Friedrichs, J. D. Chodera, R. J. Radmer, C. M. Bruns, J. P. Ku, K. A. Beauchamp, T. J. Lane, L.-P. Wang, D. Shukla, T. Tye, M. Houston, T. Stich, C. Klein, M. R. Shirts, V. S. Pande, *J. Chem. Theory Comput.* **2013**, *9*, 461.

[52] SCM, *Theoretical Chemistry*, Vrije Universiteit, Amsterdam, The Netherlands, Ams **2023**.

[53] T. P. Senftle, S. Hong, M. M. Islam, S. B. Kylasa, Y. Zheng, Y. K. Shin, C. Junkermeier, R. Engel-Herbert, M. J. Janik, H. M. Aktulga, T. Verstraelen, A. Grama, A. C. T. van Duin, *Appl Future Directions Npj Comput Mater* **2016**, *2*, 15011.

[54] J. G. Brandenburg, S. Grimme, *J. Phys. Chem. Lett.* **2014**, *5*, 1785.

[55] E. Riccardi, S. Pantano, R. Potestio, *Interface Focus* **2019**, *9*, 20190005.

[56] D. T. Zhang, L. Baldauf, S. Roet, A. Lervik, T. S. van Erp, *Proc. Natl. Acad. Sci. U. S. A.* **2024**, *121*, e2318731121.

[57] J. Behler, M. Parrinello, *Phys. Rev. Lett.* **2007**, *98*, 146401.

[58] A. P. Bartók, M. C. Payne, R. Kondor, G. Csányi, *Phys. Rev. Lett.* **2010**, *104*, 136403.

[59] S. Batzner, A. Musaelian, L. Sun, M. Geiger, J. P. Mailoa, M. Kornbluth, N. Molinari, T. E. Smidt, B. Kozinsky, *Nat. Commun.* **2022**, *13*, 2453.

[60] A. Ma, A. R. Dinner, *J. Phys. Chem. B* **2005**, *109*, 6769.

[61] H. Jung, R. Covino, A. Arjun, C. Leitold, C. Dellago, P. G. Bolhuis, G. Hummer, *Nat. Comput. Sci.* **2023**, *3*, 334.

[62] F. G. Wang, D. P. Landau, *Phys. Rev. Lett.* **2001**, *86*, 2050.

[63] A. Laio, F. L. Gervasio, *Rep. Prog. Phys.* **2008**, *71*, 126601.

## APPENDIX A: RECURSIVE BLOCK ERROR ANALYSIS

Considering Equations (1)–(3), it is easy to show that recursive block values, $k_j^r$, are identical to standard blocks averages, $k_j^s$, in case that the running estimate $k[n]$ would be a true running average (not only a 'running estimate'): a simple average of $n$ data points. However, the rate estimate $k[n]$ is not a running average in the strict sense. Instead, it is a running estimate, and it is a result obtained from different ensemble averages such as the path lengths of ensembles $[0^-]$ and $[0^+]$, and the local crossing probabilities $P_A(\lambda_{i+1}|\lambda_i)$ of ensemble $[i^+]$ with $i \geq 0$. In that case, the recursive block values and the standard block values may not necessarily be identical. In this context, the term 'the $j$th standard block value' denotes the specific value derived by segregating each relevant dataset into blocks and specifically examining the information within the $j$th block of each dataset to compute the intended property. This is denoted by the superscript $s$ in $k_j^s$, distinguishing it from a recursive block, $k_j^r$. When the recursive and standard blocks coincide (i.e., when $k_j^r = k_j^s$), we represent this as $k_j^r = k_j^s = k_j$, omitting the superscript entirely.

To discuss the recursive block approach in a more generic context, let us assume that the evaluation of $k[n]$ is based on different ensemble averages $a[n], b[n], c[n], \ldots$, that is, $k[n] = f(a[n], b[n], c[n], \ldots)$

where $f(\cdot)$ is the function that provides the rate from the ensemble averages. For these ensemble averages $a[n]$, and so forth, the running estimates are plain running averages and there is no difference between recursive and standard blocks: $a_j^r = a_j^s = a_j$ and so forth.

However, $k_j^s = f(a_j, b_j, c_j, \ldots) \neq k_j^r$ for $j > 1$. Yet, for large enough blocks we can assume that each block $j$ is close to its exact value $a, b, c$: $a_j = a + \delta a_j$, $b_j = b + \delta b_j$, $c_j = c + \delta c_j$ such that in first order of $\delta$:

$$
\begin{aligned}
k_j^s &= f(a_j, b_j, c_j, \ldots) \\
&\approx f(a, b, c, \ldots) + \left(\frac{\partial f}{\partial a}\right)\delta a_j + \left(\frac{\partial f}{\partial b}\right)\delta b_j + \left(\frac{\partial f}{\partial c}\right)\delta c_j + \cdots
\end{aligned}
\tag{A1}
$$

where the derivatives are evaluated at the exact value $a, b, c$. Now, let us compare this to the Taylor expansion of the $j$th recursive block starting from Equation (3):

$$
\begin{aligned}
k_j^r &= jf(a[jm], b[jm], c[jm], \ldots) \\
&\quad -(j-1)f(a[(j-1)m], b[(j-1)m], c[(j-1)m], \ldots) \\
&= jf\left(\frac{1}{j}\sum_{i=1}^{j} a_i, \frac{1}{j}\sum_{i=1}^{j} b_i, \frac{1}{j}\sum_{i=1}^{j} c_i, \ldots\right) \\
&\quad -(j-1)f\left(\frac{1}{j-1}\sum_{i=1}^{j-1} a_i, \frac{1}{j-1}\sum_{i=1}^{j-1} b_i, \frac{1}{j-1}\sum_{i=1}^{j-1} c_i, \ldots\right) \\
&= jf\left(a + \frac{1}{j}\sum_{i=1}^{j}\delta a_i, b + \frac{1}{j}\sum_{i=1}^{j}\delta b_i, c + \frac{1}{j}\sum_{i=1}^{j}\delta c_i, \ldots\right) \\
&\quad -(j-1)f\left(a + \frac{1}{j-1}\sum_{i=1}^{j-1}\delta a_i, b + \frac{1}{j-1}\sum_{i=1}^{j-1}\delta b_i, c + \frac{1}{j-1}\sum_{i=1}^{j-1}\delta c_i, \ldots\right)
\end{aligned}
$$

Applying Taylor expansions to $f$ in this equation and simplifying the expression, we find

$$
k_j^r \approx f(a, b, c, \ldots) + \left(\frac{\partial f}{\partial a}\right)\delta a_j + \left(\frac{\partial f}{\partial b}\right)\delta b_j + \left(\frac{\partial f}{\partial c}\right)\delta c_j + \cdots + \mathcal{O}(\delta^2)
\tag{A2}
$$

From Equations (A1) and (A2), it is evident that the Taylor expansions are equivalent up to the first order in $\delta$. This suggests that with an increase in block length, the recursive blocks converge toward the properties of standard blocks, validating our approach. It is important to realise that truncating the Taylor expansion to the first order of $\delta$ aligns with the standard error propagation practice, ensuring that our approach does not introduce any additional approximations beyond those already implied in other accepted methods.

It is interesting to note that the approach mentioned here, that is, running estimates and recursive blocks, can be effectively integrated with methods like Wang-Landau[62] or metadynamics.[63] In these techniques, a bias dynamically adapts throughout the sampling process until it converges and stabilises. Hence the bias is non-stationary and changes in the running estimate and in the running block averages. Nevertheless, as the sampling duration extends, the statistical impact of the bias' non-stationarity gradually diminishes. Therefore, when the number of MC cycles $N$ becomes sufficiently large, the effect of non-stationarity will only affect a small fraction of the blocks or only influence a minor part of the initial block. This justifies the applicability of the recursive block method even in such adaptive biasing methods.

# 11

# PAPER III (IN PREPARATION): ESTIMATING FULL PATH LENGTHS AND KINETICS FROM PARTIAL PATH TRANSITION INTERFACE SAMPLING SIMULATIONS

Manuscript in preparation.

11. Paper III (in preparation): Estimating full path
lengths and kinetics from partial path transition
interface sampling simulations

# Estimating full path lengths and kinetics from partial path transition interface sampling simulations

Wouter Vervust, Elias Wils, and An Ghysels

*IBiTech - BioMMedA group, Ghent University, The Core, Corneel Heymanslaan 10, 9000 Gent, Belgium*

Molecular dynamics (MD) simulations are crucial for investigating biological processes, yet their limitations in timescale hinder the study of rare and slow events. We recently developed replica exchange partial path transition interface sampling (REPPTIS), a path sampling method that combines the efficiency of replica exchange with the diffusive assumption of partial path ensembles, to study slow reactions with metastable states along the reactive pathways. However, REPPTIS lacks a formalism to extract time-dependent properties such as mean first passage times, fluxes, and rates. In this work, we introduce a Markov state model (MSM) framework to estimate full path lengths and kinetic properties from the overlapping partial paths generated by REPPTIS. We validate our approach using Langevin particles on one-dimensional potential energy profiles and further apply REPPTIS to the trypsin-benzamidine complex to compute dissociation kinetics. Our results highlight the challenges of sampling steep free energy landscapes with orthogonal components and suggest future improvements in the REPPTIS methodology.

## I. INTRODUCTION

Molecular dynamics (MD) simulations are essential in studying biological processes at the molecular level [1, 2]. While modern high performance computational infrastructure allows simulations to probe milliseconds, many biological processes extend largely beyond this timescale [3–5]. The challenge becomes more pronounced when kinetics are of interest, where hundreds of events are required to extract reliable statistics. However, kinetics play a pivotal role in revealing biomolecular mechanisms, where rates of conformational changes and interactions ultimately define biological function [6]. Protein-drug kinetics, for example, have been increasingly recognized for their crucial role in pharmacodynamics and better correlation with *in vivo* drug efficacy than static predictors, where continuous efforts are made to push MD simulations to longer timescales [7–10].

Path sampling methods such as replica exchange transition interface sampling (RETIS) offer an exponential speedup in extracting rate constants of rare events [11, 12]. Rare events are characterized by long waiting times (typically longer than accessible simulation times) and a rapid transit time once the event occurs. RETIS achieves this with a divide-and-conquer strategy, partitioning phase space via interfaces along an order parameter $\lambda$ between reactant state $A$ and product state $B$ (Fig. 1). Path sampling then focuses on paths that have advanced progressively further along $\lambda$, from which the rate constant can be extracted as a product of a flux term and history dependent crossing probabilities. These paths are generated by a series of Monte Carlo (MC) moves that obey detailed balance, resulting in an unbiased sampling of paths. Furthermore, paths retain all memory, meaning they are followed backward and forward in time until they reach state $A$ or $B$. These properties, full memory retention and detailed balance, protect RETIS from hysteresis effects (unfairly) favoring paths, resulting in an exact calculation of the rate constant [13–15].

The presence of long-lived metastable states along the transition pathway can, however, make RETIS paths infeasibly long, resulting in low acceptance rates and slow convergence. Partial path TIS (PPTIS) truncates path memory by confining paths to a region contained within three consecutive interfaces [16] (Fig. 1). PPTIS sacrifices exactness for computational feasibility, where the interface crossing probabilities can now be used to approximate the rate constant. While allowing a broader range of applications, PPTIS is more dependent on the choice of $\lambda$. We recently introduced the replica exchange partial path TIS (REPPTIS) methodology, which combines the TIS formalism with replica moves and partial path sampling [17]. The replica exchange move (Fig. 2C) allows paths to be extended and subsequently swapped between path ensembles, enabling paths to explore otherwise unreachable regions of phase space resulting in increased ergodic sampling.

However, both PPTIS and REPPTIS lack a formalism to extract time-dependent properties such as fluxes, rates, and mean first passage times (MFPTs), which is developed in this work. To calculate the rate constant, the conditional flux $f_A$ through $\lambda_A$ is required, which in turn depends on the length of paths crossing the interface $\lambda_A$. As PPTIS only provides *partial* paths, special effort is needed to reconstruct full paths and their lengths, incorporating an arbitrary (even infinite) number of recrossings with the intermediate interfaces.

The derivation of the flux in this paper is made possible by regarding a long MD trajectory as a Markov state model (MSM) transitioning between PPTIS ensembles, where the MSM states are consistent with the PPTIS memory assumptions. Such an analogy has been made in the past for the milestoning approach (with a very instructive explanation in Ref. [18]), where a long trajectory is seen as an MSM between subsequent milestones. As the PPTIS ensembles extend over larger and overlapping $\lambda$ regions, a new formulation is needed. After
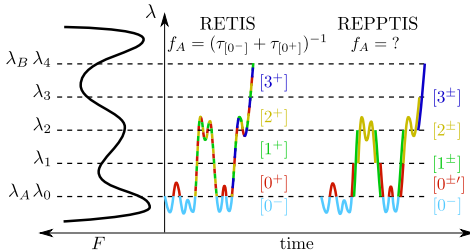
FIG. 1. Both RETIS and REPPTIS use a set of interfaces $\lambda_A$, $\lambda_1$, ..., $\lambda_B$. A long MD trajectory is cut into segments of (RE)TIS and (RE)PPTIS path ensembles. For both, sampling of state $A$ is done in the $[0^-]$ ensemble (light blue). For RETIS, paths that reach up to $\lambda_i$ are included in $[j^+]$ ($j \leq i$) ensembles. For example: red segments did not reach $\lambda_1$, and are only in $[0^+]$, while the final segment reaches $\lambda_B$, and is in all $[i^+]$ ensembles. The flux $f_A$ through $\lambda_A$ is then retrieved from the sum of average path lengths in $[0^-]$ and $[0^+]$. For REPPTIS, $[i^\pm]$ paths are confined to the interval $[\lambda_{i-1}, \lambda_{i+1}]$. An exception is $[0^{\pm\prime}]$, with paths confined to the $[\lambda_A, \lambda_1]$. The long MD segment is now decomposed into overlapping path segments, for which average full path length statistics will be retrieved in this work.

deriving the MSM, it is used to compute the flux and other mean first passage times.

The paper starts with a short review of notations, ensemble definitions, and global and local crossing probabilities. Next, we introduce the view of a long MD trajectory as a sequence of overlapping PPTIS segments, which leads to the Markov state model with transition probability matrix $M$. Using the MSM, we show how the average path lengths can be retrieved, which gives us MFPTs under various boundary conditions. The global crossing probability also emerges from the MSM, giving us a closed-form formula for $P_A(\lambda_A \to \lambda_B)$, as an alternative to the known iterative procedure. The MSM formalism is then applied to several systems to validate the new equations for the flux and MFPTs. REPPTIS is then applied to study the dissociation kinetics of the trypsin-benzamidine complex, after which the results are discussed and conclusions are drawn.

## II. BUILDING MSM FOR REPPTIS

First, the notations of TIS and PPTIS are reviewed. TIS will serve as reference for validation of the new PP-TIS flux equations. Next, the PPTIS method and its $[i^\pm]$ path ensembles are described in more detail, where the memory assumptions of PPTIS are used to interpret a long MD trajectory as a Markov state model (MSM) that jumps between the PPTIS path ensembles. It is then discussed how properties, such as long crossing probabilities and mean first passage times (MFPTs) can be obtained

from this specific MSM. This also leads to the equation giving the PPTIS flux.

### A. Notations for TIS and PPTIS

In both TIS and PPTIS, a set of non-intersecting interfaces $\lambda_0 = \lambda_A$, $\lambda_1$, ..., $\lambda_{N-1} = \lambda_B$ are distributed along an order parameter $\lambda$. The first and last interfaces define the reactant state $A$ ($\lambda < \lambda_A$) and product state $B$ ($\lambda > \lambda_B$). Path ensembles $[i^+]$ (TIS) and $[i^\pm]$ (PPTIS) are associated to these interfaces (Figs. 2A-B), which are sampled using an MC approach in path space (shooting move, replica exchange move, etc.). TIS, designed for rare events, allows an exact calculation of the forward rate $k_{AB}$

$$k_{AB} = f_A P_A(\lambda_A \to \lambda_B), \qquad (1)$$

where $f_A$ is the conditional flux (leaving) through $\lambda_A$, and $P_A(\lambda_A \to \lambda_B)$ is the global or overall crossing probability, i.e. the probability that a path that has just left state $A$ will cross $\lambda_B$ before recrossing with $\lambda_A$. PPTIS, designed for events that are both rare as well as slow, sacrifices exactness by truncating path memory, resulting in an approximation of the global crossing probability. In-depth discussions on TIS ensembles can be found in Refs. [14, 15], whereas PPTIS ensembles are covered in more detail below.
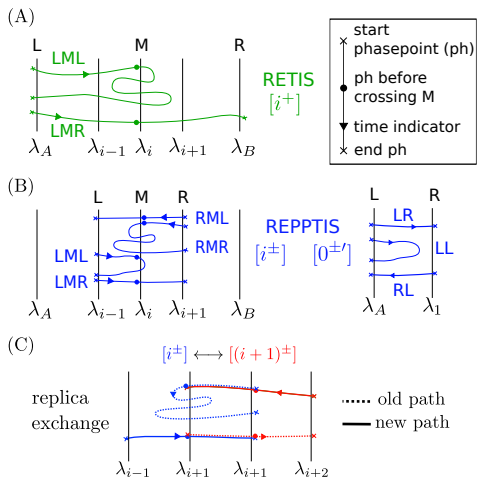


FIG. 2. **A**: Paths of the TIS $[i^+]$ ensemble **B**: Paths of the PPTIS $[i^\pm]$ ensemble and the special case $[0^{\pm\prime}]$ ensemble. **C**: In the replica exchange move of REPPTIS, paths of adjacent ensembles are extended and subsequently exchanged.

PPTIS path ensembles $[i^\pm]$ ($\forall i \in [1, N-2]$) contain all paths that cross $\lambda_i$ (middle, label M), and start and

11. Paper III (in preparation): Estimating full path
lengths and kinetics from partial path transition
interface sampling simulations

end in the neighboring interfaces $\lambda_{i-1}$ (left, label L) and $\lambda_{i+1}$ (right, label R), as visualized in Fig. 2B. Thus, each ensemble contains four path types (LML, LMR, RML, and RMR). For example in $[i^\pm]$, an LMR path starts at $\lambda_{i-1}$ (L), crosses $\lambda_i$ (M) before recrossing $\lambda_{i-1}$, and then ends at $\lambda_{i+1}$ (R) before recrossing with $\lambda_{i-1}$.

Near the $\lambda_0$ interface, two additional path ensembles $[0^-]$ and $[0^{\pm\prime}]$ are defined. The ensemble $[0^-]$ is used for both TIS and PPTIS simulations. It contains all paths that start at $\lambda_0$, travel into the reactant state $\lambda < \lambda_0$, and end at $\lambda_0$. One can define $\lambda_{0^-} = \lambda_0 - \delta$ as the M interface for $[0^-]$, where $\delta$ is an infinitesimal positive number. As all paths in $[0^-]$ automatically cross this interface, its presence remains implicit, and $[0^-]$ paths are denoted as RR path types (instead of RMR). The ensemble $[0^{\pm\prime}]$ (Fig. 2B) contains all paths that (1) start at $\lambda_0$ and end at either $\lambda_0$ or $\lambda_1$, or (2) start at $\lambda_1$ and end at $\lambda_0$. One could similarly define $\lambda_{0^+} = \lambda_0 + \delta$ as the M interface of $[0^{\pm\prime}]$, where $\delta$ is an infinitesimal positive number. Also here, the presence of this interface is implied, resulting in LR, LL and RL path type notations. RR type paths in $[0^{\pm\prime}]$ have an infinitesimal weight as the M interface is practically equivalent to $\lambda_0$, and they are not required for any of the analysis that follows.

Local crossing probabilities $p_{[i^\pm]}^{k,l}$ ($\forall i \in [1, N-2]$, $k, l \in [-1, +1]$) can be calculated from the PPTIS ensembles $[i^\pm]$. Here, $p_{[i^\pm]}^{k,l}$ denotes the probability that a path that crossed $\lambda_i$ right after $\lambda_{i+k}$ will cross $\lambda_{i+l}$ before crossing $\lambda_{i-l}$. For example, the local crossing probability $p_{[i^\pm]}^{-1,-1}$ of $[i^\pm]$ denotes the probability that a path that crossed $\lambda_i$ right after $\lambda_{i-1}$ will cross $\lambda_{i-1}$ before crossing $\lambda_{i+1}$. This is estimated from the simulation output as the ratio of $\text{LML}_{[i^\pm]}$ paths to the total number of $\text{LML}_{[i^\pm]}$ and $\text{LMR}_{[i^\pm]}$ paths. These local crossing probabilities can then be used to estimate the global crossing probability $P_A(\lambda_A \to \lambda_B)$.

The PPTIS ensembles are sampled using a MC scheme, where newly generated paths are subjected to a Metropolis acceptance rule [19], enforcing detailed balance. Having microscopic reversibility has the advantage that paths are sampled according to their respective ensemble weights, as the non-localized sampling ensures ergodicity. Large barriers orthogonal to the chosen $\lambda$ parameter can, however, highly restrict ergodic sampling. As such, the construction of specialized MC moves that facilitate path decorrelation and increase phase space exploration is of high importance. Inclusion of a replica exchange move greatly increased ergodic sampling for both TIS and PPTIS [11, 17]. Other MC moves include stone skipping, web throwing, wire-fencing, target-swap move, mirror move, and more [17, 20–22]. A downside of enforcing detailed balance is a limitation to equilibrium situations, unlike splitting based methods such as forward flux sampling [23].

## B. Building the Markov state model

Consider a very long equilibrium trajectory, visualized in Fig. 3A. The trajectory can be decomposed into (overlapping) path segments that are part of the PPTIS path ensembles. These are the different colored path segments of Fig. 3. The long MD trajectory $U$ is the sequence of (overlapping) path segments $(U_0, U_1, U_2, \ldots)$, where each path segment $U_n$ is part of a specific path type of a PPTIS path ensemble. For the trajectory in Fig. 3, $U_0 \in \text{LR}_{[0^{\pm\prime}]}$, $U_1 \in \text{LMR}_{[1^\pm]}$, $U_2 \in \text{LML}_{[2^\pm]}$, $U_3 \in \text{RML}_{[1^\pm]}$, etc.

The four path types of $[i^\pm]$ can be viewed as four states $S_i^{k,l}$, where $k$ and $l$ denote the starting and ending interfaces of the path type, respectively. The values $k, l$ lie in $\{-1, +1\}$ with $-1$ referring to L and $+1$ referring to R. Thus, $S_i^{-1,-1}$ corresponds to $\text{LML}_{[i^\pm]}$, the LML paths in $[i^\pm]$. Similarly, $S_i^{-1,+1}$ refers to $\text{LMR}_{[i^\pm]}$, $S_i^{+1,-1}$ to $\text{RML}_{[i^\pm]}$, and $S_i^{+1,+1}$ to $\text{RMR}_{[i^\pm]}$. This new notation simplifies the equations later on (Eq. 6). The long trajectory $U = (U_0, U_1, U_2, \ldots)$ of Fig. 3A can thus be mapped to the state chain $\left( S_0^{-1,+1}, S_1^{-1,+1}, S_2^{-1,-1}, S_1^{+1,-1}, \ldots \right)$. The state space $\mathcal{S}$ is defined as the set of all states $S_i^{k,l}$, where $i = 0, \ldots, N-1$ and $k, l \in \{-1, +1\}$. The limiting states related to states $S_i^{k,l}$ with $i$ equal to 0, $N-2$, or $N-1$ are specific and are discussed in subsection II C. For all other $i$, there are four states associated to a path ensemble $[i^\pm]$.

Looking at all the segments in the long trajectory $U$, some statistics can be done. For instance, let us select all the segments that are in state $S_1^{-1,+1}$, i.e. they are an LMR path in $[1^\pm]$ connecting $\lambda_0 \to \lambda_1 \to \lambda_2$. For each of those selected segments, the next interface that possibly can be hit by the long trajectory must be either a return to $\lambda_1$ or a progression to $\lambda_3$ (see Fig. 3). This would result in an LML path (in case of return) or LMR path (in case of progression) in $[2^\pm]$, corresponding to the states $S_2^{-1,-1}$ or $S_2^{-1,+1}$, respectively. In a long trajectory $U$, a fraction of the transits will be to $S_2^{-1,-1}$ and the other fraction to $S_2^{-1,+1}$. The probability of the long trajectory going from $S_1^{-1,+1}$ to $S_2^{-1,+1}$ is the probability that a path that has just crossed $\lambda_2$ after having crossed $\lambda_1$ will cross $\lambda_3$ before recrossing $\lambda_1$. This is exactly the PPTIS local crossing probability $p_{[2^\pm]}^{-1,+1}$ which can be computed directly from counting the paths in the $[2^\pm]$ ensemble. The other fraction is then $p_{[2^\pm]}^{-1,-1} = 1 - p_{[2^\pm]}^{-1,+1}$.

In this example, the transition probabilities are given by

$$P\left( S_1^{-1,+1} \to S_2^{-1,+1} \right) = p_{[2^\pm]}^{-1,+1} \tag{2}$$

$$P\left( S_1^{-1,+1} \to S_2^{-1,-1} \right) = p_{[2^\pm]}^{-1,-1} \tag{3}$$

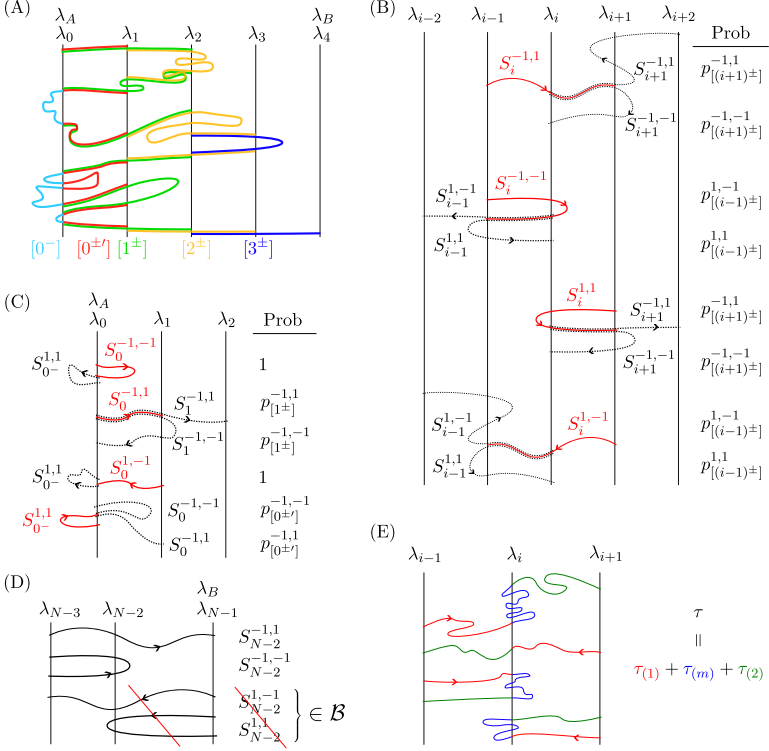and the transition probability to any of the other states

FIG. 3. **A**: A long MD path decomposed into its (overlapping) PPTIS path segments. **B**: Ensembles $[i^{\pm}]$ $(i \in [2, N-3])$ contain 4 path types (red) that can transition to two possible path types in a neighboring ensemble (black dotted) with a probability equal to the PPTIS local crossing probabilities. **C**: The ensembles near $\lambda_A$. The $[0^-]$ ensemble contains only one state, and the $[0^{\pm\prime}]$ ensemble contains three states. All possible transitions and their probabilities are shown. **D**: The ensemble $[(N-2)^{\pm}]$ near $\lambda_B$. The two states denoting paths that start from $\lambda_B$ are combined into the single state $S_B$. **E**: Path lengths $\tau$ are decomposed into three parts: the part $\tau_{(1)}$ before the first crossing of $\lambda_i$, the part $\tau_{(2)}$ after the last crossing of $\lambda_i$, and the part $\tau_{(m)}$ in between. These parts can be zero, as visible in the second path (no middle part).

is zero,

$$P\left(S_1^{-1,+1} \to S_i^{k,l}\right) = 0 \qquad (4)$$

More generally, the transition probabilities can be gathered in the transition matrix $M$,

$$M_{ikl,i'k'l'} = P\left(S_i^{k,l} \to S_{i'}^{k',l'}\right) \qquad (5)$$

with elements

$$M_{ikl,i'k'l'} = \begin{cases} p_{i+l}^{-l,+1}, & i' = i+l, k' = -l, l' = +1 \\ p_{i+l}^{-l,-1}, & i' = i+l, k' = -l, l' = -1 \\ 0, & \text{elsewhere} \end{cases} \qquad (6)$$

While the indices look tedious, it means that the end point $l$ of the segment $S_i^{k,l}$ determines whether the next state is in a higher ensemble $i' = i + 1$ ($l = +1$, so end point R) or in a lower ensemble $i' = i - 1$ ($l = -1$, so end point L). This explains that the next visited path ensemble is surely $[i'^{\pm}]$ with $i' = i+l$. Moreover, the end point $l$ determines the starting point of the next segment, so the next starting point is $k' = -l$. The new segment can then have an arbitrary end point, so $l'$ is either $+1$ or -1.

Interestingly, the elements of $M_{ikl,i'k'l'}$ are independent of the starting point of the $S_i^{k,l}$ segment, see Eq. 6. This implies that the memory about the initial starting

point $k$ is lost. This is indeed implied by the PPTIS ensembles $[i^\pm]$ which cover segments with three interface labels but not four labels. Therefore an initial path with labels $(i, k, l)$ will lose information on the label $k$ when a new state is acquired. (The labels $i$ and $l$ survive through the updated labels $i' = i + l$ and $k' = -l$.) Hence, the proposed transition matrix indeed restrains the memory to three labels, in accordance with the PPTIS ensembles.

The long MD trajectory can thus be translated to a Markov chain between the states $S_i^{k,l}$ with transition probabilities given by the matrix elements of $M$, assuming the memory covers three labels. For convenience, we introduce an alternative notation where the three labels $ikl$ are flattened to a Greek index $\alpha$. The state space $\mathcal{S}$ can then be written as $(S_1, S_2, \ldots, S_\alpha, \ldots)$ with transition probabilities $M_{\alpha\beta} = P(S_\alpha \to S_\beta)$.

Prior PPTIS work has consistently used $p_i^=$, $p_i^\pm$, $p_i^\mp$, and $p_i^\ddagger$ to denote the local crossing probabilities associated to $\mathrm{LML}_{[i^\pm]}$, $\mathrm{LMR}_{[i^\pm]}$, $\mathrm{RML}_{[i^\pm]}$, and $\mathrm{RMR}_{[i^\pm]}$, respectively. Here, the $p_{[i^\pm]}^{k,l}$ notation is introduced to enable general equations such as Eq. 6. All possible transitions of $[i^\pm]$ paths are shown in Fig. 3B, together with their state notation and associated transition probabilities.

### C. Limiting states in MSM

For all path ensembles $i = 1, \ldots, N - 2$, there are four states $S_i^{k,l}$ with $k, l \in \{+1, -1\}$. The number of states is different for the limiting ensembles associated to $\lambda_0$, $\lambda_{N-2}$, and $\lambda_{N-1}$. The corresponding states and transition probabilities are shown in Figs. 3C-D, and as a graph representation focusing on the limiting ensembles in Fig. 4.

For the paths in path ensemble $[0^-]$, there are only $k = l = +1$ paths, so this ensemble contributes only one state denoted as state $S_{0^-}^{+1,+1}$. The paths in path ensemble $[0^{\pm\prime}]$ can be of type LL, LR, or RL, and contribute three states (Fig. 3C).

Next, consider the paths in $\mathcal{B}$, which is the path ensemble of all paths where the last visited region was $B$ rather than $A$. All paths start at $\lambda_{N-1}$ and eventually return to region $A$ with probability 1 in a very long equilibrated MD simulation. These segments are grouped into state $S_{\mathcal{B}}$.

The paths in path ensemble $[(N-2)^\pm]$ can in principle also be of all four path types LML etc. However, an RMR or RML segment in $[(N-2)^\pm]$ started by crossing $\lambda_{N-1}$, which implies that the trajectory was last in region $B$. Such segments are thus part of $\mathcal{B}$, and the states $S_{N-2}^{+1,+1}$ or $S_{N-2}^{+1,-1}$ are thus not included as separate states in state space $\mathcal{S}$ (Fig. 3D).

This brings the number of states in $\mathcal{S}$ to $4N - 5$: 1 for $[0^-]$, 3 for $[0^{\pm\prime}]$, 2 for $[(N-2)^\pm]$, 1 for $\mathcal{B}$, and 4 for each of the other $N - 3$ ensembles.

As an example, consider a PPTIS simulation with



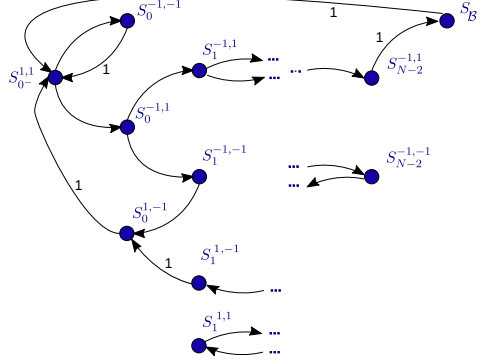FIG. 4. Visualization of the connectivity in the MSM network for the limiting ensembles around state $A$ and $\mathcal{B}$. Each circle represents a state. Transitions for $[i^+]$ states ($i \in [2, N - 3]$) are shown in Fig. 3B. Transitions with probability 1 are indicated.

$N = 5$ interfaces ($\lambda_0$, $\lambda_1$, $\lambda_2$, $\lambda_3$, $\lambda_4$) with $4N - 5 = 15$ states. The $15 \times 15$ transition matrix $M$ is given in Eq. 7 in Table I, where shorthand notation $p_i^{k,l}$ is used to denote $p_{[i^\pm]}^{k,l}$. The states are ordered first by $i$, then by $k$, and finally by $l$. The sum of each row is equal to 1. The elements that are not shown are equal to zero. The elements $p_0^{-1,+1}$ and $p_0^{-1,-1}$ were denoted as $p_0^\pm$ and $1 - p_0^\pm$ in the REPPTIS paper [17]. A certain block matrix structure can be recognized, because states related to $\lambda_i$ only connect with neighboring states related to $\lambda_{i+1}$ or $\lambda_{i-1}$. This gives zero blocks along the matrix diagonal.

Since the connection has now been completely made between the long MD trajectory and a Markov state model governed by $M$, several concepts from Markov state modeling can be applied, such as the propagation of an initial state distribution, the stationary distribution, hitting probability, and recurrent probability.

### D. Global crossing probability from MSM

For the rate calculation, the probability of interest is the global crossing probability $P_A(\lambda_A \to \lambda_B)$, which is the probability to cross $\lambda_B$ before returning crossing $\lambda_A$, given that $\lambda_A$ is crossed at this moment. The MSM for the REPPTIS ensembles can be used to compute this type of probability. In general MSM theory, it is well known how to derive so-called hitting and return probabilities from the matrix $M$ under various boundary conditions (see SI). For the specific case of $P_A(\lambda_A \to \lambda_B)$, two such concepts are needed: 1) probability $Z$ to reach state $\alpha$ before $\beta$, starting from $\delta$, and 2) probability $Y$ to reach state $\alpha$ before $\beta$, starting from $\alpha$, *and* leaving $\alpha$ in the first step.

TABLE I. Example of the $15 \times 15$ matrix $M$ for 5 interfaces ($\lambda_0$, $\lambda_1$, $\lambda_2$, $\lambda_3$, $\lambda_4$).

| FROM\TO | $S_{0-}$ | $S_0^{-1,-1}$ | $S_0^{+1,-1}$ | $S_0^{+1,-1}$ | $S_1^{-1,-1}$ | $S_1^{+1,-1}$ | $S_1^{+1,-1}$ | $S_1^{+1,+1}$ | $S_2^{-1,-1}$ | $S_2^{-1,+1}$ | $S_2^{+1,-1}$ | $S_2^{+1,-1}$ | $S_3^{-1,-1}$ | $S_3^{-1,+1}$ | $S_B$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $S_{0-}^{+1,+1}=S_A$ | | $p_0^{-1,-1}$ | $p_0^{-1,+1}$ | 0 | | | | | | | | | | | |
| $S_0^{-1,-1}$ | 1 | | | | 0 | 0 | 0 | 0 | | | | | | | |
| $S_0^{-1,+1}$ | 0 | | | | $p_1^{-1,-1}$ | $p_1^{-1,+1}$ | 0 | 0 | | | | | | | |
| $S_0^{+1,-1}$ | 1 | | | | 0 | 0 | 0 | 0 | | | | | | | |
| $S_1^{-1,-1}$ | | **1** | | | | | | | 0 | 0 | 0 | 0 | | | |
| $S_1^{-1,+1}$ | | 0 | | | | | | | $p_2^{-1,-1}$ | $p_2^{-1,+1}$ | 0 | 0 | | | |
| $S_1^{+1,-1}$ | | **1** | | | | | | | 0 | 0 | 0 | 0 | | | |
| $S_1^{+1,+1}$ | | 0 | | | | | | | $p_2^{-1,-1}$ | $p_2^{-1,+1}$ | 0 | 0 | | | |
| $S_2^{-1,-1}$ | | | | | 0 | 0 | $p_1^{+1,-1}$ | $p_1^{+1,+1}$ | | | | | 0 | 0 | |
| $S_2^{-1,+1}$ | | | | | 0 | 0 | 0 | 0 | | | | | $p_3^{-1,-1}$ | $p_3^{-1,+1}$ | |
| $S_2^{+1,-1}$ | | | | | 0 | 0 | $p_1^{+1,-1}$ | $p_1^{+1,+1}$ | | | | | 0 | 0 | |
| $S_2^{+1,+1}$ | | | | | 0 | 0 | 0 | 0 | | | | | $p_3^{-1,-1}$ | $p_3^{-1,+1}$ | |
| $S_3^{-1,-1}$ | | | | | | | | | 0 | 0 | $p_2^{+1,-1}$ | $p_2^{+1,-1}$ | | | 0 |
| $S_3^{+1,+1}$ | | | | | | | | | 0 | 0 | 0 | 0 | | | 1 |
| $S_B$ | 1 | | | | | | | | | | | | 0 | 0 | |

$$M = \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad (7)$$

We first recap the equations for $Z$ and $Y$ with general states (see SI for a full derivation) and then apply them to compute the desired $P_A(\lambda_A \to \lambda_B)$.

Define the $Z_{\delta(\beta)}$ as the probability that a trajectory starting in state $S_\delta$ can reach state $S_\beta$ before reaching state $S_\alpha$,

$$Z_{\delta(\beta)} \equiv P(\exists n \ge 0 : S_n = \beta \text{ before } S_n = \alpha | S_0 = \delta) \quad (8)$$

which is a hitting probability. The equations for $Z$ are drafted by considering different starting values $\delta$, where the states $\alpha, \beta$ are fixed in this paragraph. Starting from state $\delta = \beta$, you are already in the desired state with certainty, so the probability is 1. However, starting from $\delta = \alpha$, it is certain that you will not reach $\beta$ first, because you already reached the $\alpha$ boundary. This gives the following $Z$ values for these boundaries,

$$Z_{\beta(\beta)} = 1 \qquad (9)$$
$$Z_{\alpha(\beta)} = 0 \qquad (10)$$

When starting from $\delta \ne \beta$ and $\delta \ne \alpha$, you need to take at least one step ($n \ge 1$). By conditioning the probability on this first step, the other equations become (see SI)

$$Z_{\delta(\beta)} = \sum_\gamma M_{\delta\gamma} Z_{\gamma(\beta)} \quad \delta \ne \alpha, \beta \qquad (11)$$

Eqs. 9-11 are sufficient to determine all $Z_{\delta(\beta)}$, $\forall \delta$.

Next, we compute the probability $Y$ that a trajectory that left state $S_\alpha$, can reach state $S_\beta$ before returning to the initial state $S_\alpha$,

$$Y_{\alpha(\beta)} \equiv P(\exists n \ge 1 : S_n = \beta \text{ before } S_n = \alpha | S_0 = \alpha) \quad (12)$$

This probability is the complement of the *return* probability to state $\alpha$. It is assumed that $\alpha \ne \beta$. The number of steps must be at least $n = 1$. By conditioning on the first step, the probability to reach $\beta$ before returning to

$\alpha$ can be written in terms of the previous $Z$ vectors (see SI),

$$Y_{\alpha(\beta)} = \sum_\gamma M_{\alpha\gamma} Z_{\gamma(\beta)} \qquad (13)$$

Eq. 13 can be used to calculate $Y_{\alpha(\beta)}$ straightforwardly from $M$ and the previously computed $Z_{\delta(\beta)}$ components.

Finally, the global crossing probability for REPPTIS can be computed by making a specific choice for $\alpha$, $\beta$, and $\delta$ in the calculation of $Z$ and $Y$. The global crossing probability $P_A(\lambda_A \to \lambda_B)$ is the probability to cross $\lambda_B$ before crossing $\lambda_A$, given that $\lambda_A$ is crossed at this moment. This corresponds to a probability of the type $Y_{\alpha(\beta)}$ as in Eq. 12. It can be obtained by setting $\alpha = S_{0-}^{+1,+1}$ in the initial reactant state and $\beta = S_B$ in the product state. The first step is to compute the $Z$ vector (see SI for solving Eqs. 9-11), followed by the computation of the $Y_{\alpha(\beta)}$ component with Eq. 13,

$$P_A(A \to B) = Y_{\alpha(\beta)} = \sum_\gamma M_{\alpha\gamma} Z_{\gamma(\beta)} \qquad (14)$$

We have derived a new MSM-based result for the global crossing probability, which is equivalent to solving the set of recursive (RE)PPTIS equations in previous publications [16, 17, 24, 25]. The equivalence of both methods can intuitively be understood as follows. The solution to Eqs. 9-11 gives the MSM-based result, and it requires the matrix inverse of an adapted form of $M$ (see SI). Such matrix inverse implies continued fractions of its elements, represented by the local crossing probabilities $p_{[i\pm]}^{k,l}$. Use of such continued fractions is also inherent to the iterative scheme used in the previous (RE)PPTIS equations, hinting to equivalence. Moreover, it was numerically verified that the MSM-based results are identical to the results of the iterative scheme.

## III.   KINETICS WITH MSM FOR REPPTIS

In the previous section, it was demonstrated how a Markov state model can be built for REPPTIS and how the crossing probability can be extracted. This section moves from probabilities to kinetic quantities, where the path lengths of the PPTIS segments $U_n$ also come into play. Average path lengths have a unit of time, and thus give information about the relevant timescales of the studied biological or chemical process. Examples of interesting path lengths are the averages $\tau_{[0^-]}$ and $\tau_{[0^+]}$, which are the average time spent per path in the $[0^-]$ and $[0^+]$, respectively. Their sum leads to the conditional flux $f_A = \left(\tau_{[0^-]} + \tau_{[0^+]}\right)^{-1}$ and subsequently to the reaction rate $k = f_A\, P_A(A \to B)$.

The average path length $\tau_{[0^+]}$ is however not directly available from the REPPTIS output, as the PPTIS path ensembles consist of segments rather than full paths. Consequently, the MSM formalism will be used to compute average path lengths for general paths that extend beyond a single MSM state. Our approach differs from the standard calculation of hitting times in MSM networks, which focuses on the number of jumps. Here, the focus lies on the amount of accumulated time. Roughly speaking, the essence of the method is that walks through the MSM network are built randomly. The walker jumps from state to state with a certain hopping probability (the transition probabilities in $M$). At every jump to a new state, the walker accumulates more time in the new state, until the walker reaches its destination. By stitching together the segments, and only counting time of the non-overlapping parts of the segments, the average path length is obtained.

This type of average path length is also commonly referred to as a mean first passage time (MFPT), where the first passage refers to reaching the destination of the walker. In the following subsections, notations to distinguish the non-overlapping parts will be introduced, the general equations for MFPTs of an MSM will be reviewed with attention for the different conditioning choices (e.g. destination criteria), and the equations will be applied to compute the path length $\tau_{[0^+]}$ and the flux.

### A.   Path lengths and overlap

Introduce the average path length of a state $S_i^{k,l}$ as $\tau_i^{k,l}$. In the flattened notation, this reads as an average path length $\tau_\alpha$ or state $\alpha$. Concretely, $\tau_\alpha$ is computed by counting the number of phase points in each of the paths and by averaging these lengths over all the paths in $S_\alpha$.

In order to measure the total time in an MD trajectory, one should avoid double counting the time spent in the overlapping parts of some of the segments. As shown in Fig. 3E, each path in a state $S_i^{k,l}$ can be divided in three pieces: the first part (1) *before* crossing $\lambda_i$ for the first time, the last part (2) *after* crossing $\lambda_i$ for the last time,

| | | | $\lambda_i$ | (1) | (m) | (2) |
|---|---|---|---|---|---|---|
| $[0^-]$ | RR | $S_0^{+1,+1}$ | $\lambda_0$ | 0 | x | 0 |
| $[0^\pm]$ | LL | $S_0^{-1,-1}$ | $\lambda_0$ | 0 | x | 0 |
| $[0^\pm]$ | LR | $S_0^{-1,+1}$ | $\lambda_0$ | 0 | 0 | x |
| $[0^\pm]$ | RL | $S_0^{+1,-1}$ | $\lambda_0$ | x | 0 | 0 |
| $[i^\pm]$ | *M* | $S_i^{k,l}$ | $\lambda_i$ | x | x | x |
| $\mathcal{B}$ | RM* | $S_{\mathcal{B}}$ | $\lambda_B$ | 0 | x | x |

TABLE II. The three contributions to path length for special path segments: first part (1) *before* crossing $\lambda_i$ for the first time, last part (2) *after* crossing $\lambda_i$ for the last time, and middle part (m) in *between*. The star sign $*$ can be either L or R.

and the middle part (m) *between* crossing $\lambda_i$ for the first and last time. The average path length in $S_\alpha$ can thus be written as the sum of three averages,

$$\tau_\alpha = \tau_{(1),\alpha} + \tau_{(m),\alpha} + \tau_{(2),\alpha} \tag{15}$$

This is generally true even if some of the paths in state $S_i^{k,l}$ might lack a middle piece. The cases for $[0^-]$ and $[0^{\pm\prime}]$ are quite specific, where we make use of the implied (m)iddle interface $\lambda_{0-} = \lambda_0 - \delta$ for $[0^-]$ and $\lambda_{0+} = \lambda_0 + \delta$ for $[0^{\pm\prime}]$, with $\delta$ an infinitesimal positive number. In $[0^-]$, part (1) and part (2) are then zero in each path, such that the average over part (m) is actually equal to the average path length of the $[0^-]$ ensemble.

$$\tau_{[0^-]} = \tau_{(m),[0^-]} \tag{16}$$

In $[0^{\pm\prime}]$, each path type has only one nonzero path,

$$\tau_{[0^{\pm\prime}]}^{-1,-1} = \tau_{(m),[0^{\pm\prime}]}^{-1,-1} \tag{17}$$

$$\tau_{[0^{\pm\prime}]}^{-1,+1} = \tau_{(2),[0^{\pm\prime}]}^{-1,+1} \tag{18}$$

$$\tau_{[0^{\pm\prime}]}^{+1,-1} = \tau_{(1),[0^{\pm\prime}]}^{+1,-1} \tag{19}$$

Without loss of generality, the paths can thus still be split up in the three parts with Eq. 15, even if some parts are zero. This is summarized in Table II. In practice, the three different parts in the average path length can be detected in every state $S_i^{k,l}$, by detecting the first and last crossing points with $\lambda_i$ in each path of $S_i^{k,l}$ that is sampled by (RE)PPTIS. Going from state $S_i^{k,l}$ to state $S_{i+l}^{-l,l'}$ can now be seen as accumulation of the extra time $\tau_{(m),i+l}^{-l,l'} + \tau_{(2),i+l}^{-l,l'}$ to the full path length, where the first part (1) is skipped to avoid double counting of the overlapping path segments.

### B.   General MFPT equations

In this subsection, the general equations for mean first passage times are drafted. Similarly to section II D, let us call $\delta$ a general starting state for the path. The final destinations of the paths are collected in a set $C$. If a path reached any of the possible destinations in $C$ for the

first time after $n$ steps, the accumulated time is denoted as $t_n^C$. The equations for two types of MFPTs will be computed: 1) the average time $G$ to reach state $\alpha$ or $\beta$ starting from $\delta$, and 2) the average time $H$ to reach state $\alpha$ or $\beta$ starting from $\alpha$ and leaving $\alpha$ in the first step (a return time).

Assume $G_{\delta(C)}$ is the average time to reach state $\alpha$ or $\beta$ starting from state $\delta$. This gives $C = \{\alpha, \beta\}$ as the set of destinations for the paths. On average, this gives the quantity $G_{\delta(C)}$,

$$G_{\delta(C)} \equiv E(t_n^C, n \geq 0 | S_0 = \delta) \tag{20}$$

For $\delta \in C$, no steps need to be taken, and no time is accumulated,

$$G_{\delta(C)} = 0, \forall \delta \in C \tag{21}$$

For $\delta \notin C$, at least one step is taken, e.g. from state $\delta$ to state $\gamma$ with probability $M_{\delta\gamma}$. By conditioning on this first step, a recursive relation can be built,

$$G_{\delta(C)} = \tau_{(m2),\delta} + \sum_{\gamma} M_{\delta\gamma} G_{\gamma(C)} \tag{22}$$

where the first term $\tau_{(m2),\delta} = \tau_{(m),\delta} + \tau_{(2),\delta}$ refers to the time that is accumulated in this step. In the second term, Eq. 22 can be re-applied on $G_{\gamma(C)}$, and so on, until the destination is reached and Eq. 21 stops the time accumulation.

Next, assume $H_{\delta(C)}$ is the average time to reach state $\alpha$ or $\beta$ starting from state $\delta$, given that at least one step has been taken to leave $\alpha$. This gives $C = \{\alpha, \beta\}$ as the set of destinations. The times $H_{\delta(C)}$ are identical to $G_{\delta(C)}$ for all $\delta \neq \alpha$. Nevertheless, the case of a starting position $\delta = \alpha$ is the one of interest, which gives $H$ the interpretation of a return time. Conditioning on the obligatory first step gives an expected stopping time $H_{\alpha(C)}$,

$$H_{\alpha(C)} \equiv E(t_n^C, n \geq 1 | S_0 = \alpha) \tag{23}$$

$$= \tau_{(m2),\alpha} + \sum_{\gamma} M_{\alpha\gamma} G_{\gamma(C)} \tag{24}$$

By first computing the $G_{\delta(C)}$ vector with Eqs. 21-22, $H_{\alpha(C)}$ can be computed as well with Eq. 24.

As a last note, the contribution of the starting state $\delta$ might need to be modified depending on the application of these mean first passage times $G$ and $H$. E.g. it can be reduced to only include part (2) instead of parts (m2), which can be achieved by subtracting $\tau_{(m),\delta}$. This can be needed to impose that the counting starts after the last crossing of $\lambda_i$. Another modification can be that the starting state $\delta$ does not contribute at all, which is achieved by extracting $\tau_{(m2),\delta}$. Such a modification moreover only makes sense if the path length is at least one step.

## C. Application of MFPT equations

Specifically for the flux $f_A$, the average path lengths in the $[0^-]$ and $[0^+]$ ensembles should be computed and summed to obtain $f_A = \left(\tau_{[0^-]} + \tau_{[0^+]}\right)^{-1}$. Here, $\tau_{[0^-]} \equiv \tau_{(m),0^-}$ which is directly obtained from the REPPTIS simulation output. An estimation for $\tau_{[0^+]}$ is not directly available from the output. It can be computed as the mean first passage time of leaving $\alpha = S_{0^-}^{+1,+1}$ and either returning to $\alpha = S_{0^-}^{+1,+1}$ or reaching $\beta = S_{\mathcal{B}}$. This gives a starting position $\alpha = S_{0^-}^{+1,+1}$ while $C = \{S_{0^-}^{+1,+1}, S_{\mathcal{B}}\}$ is the set of possible destinations. The path should leave the initial state, so at least one step must be taken, leading to

$$\tau_{[0^+]} = H_{\alpha(C)} - \tau_{(m),\alpha} \tag{25}$$

Here, the (m) part of the first state was subtracted to start counting time when $\lambda_0$ is last crossed.

As an alternative quantity, the sum $\tau_{[0^-]} + \tau_{[0^+]}$ can be directly computed. It represents the average time of a path measured from the moment it first entered region A until the moment the path enters region B. It can be computed as the mean first passage time of leaving $\alpha = S_{0^-}^{+1,+1}$ and reaching $\beta = S_{\mathcal{B}}$, possibly after multiple revisits to state $\alpha = S_{0^-}^{+1,+1}$. Compared to the previous equation, this gives the same starting position $\alpha = S_{0^-}^{+1,+1}$. However, the set of destinations is now $C = \{S_{\mathcal{B}}\}$, and unlike in the previous equation, the first state does not need to be adapted, giving

$$\tau_{[0^-]} + \tau_{[0^+]} = G_{\alpha(C)} \tag{26}$$

The interpretation of this time is the time between entering and leaving region A. Therefore, on average, one exit from region A will be registered reaching region B in a time span $\tau_{[0^-]} + \tau_{[0^+]}$, indeed leading to the flux $f_A = 1/(\tau_{[0^-]} + \tau_{[0^+]})$. In the examples of Sec. IV we have verified that the flux derived from Eq. 25 is identical to the flux derived from Eq. 26

## IV. ILLUSTRATION WITH 1D SYSTEMS

The MSM framework is applied to REPPTIS simulations of a one-dimensional particle diffusing in several (smooth) potential landscapes $V(x)$, where $x$ denotes the position of the particle. The analytical expressions for the considered potentials are given in the Supplementary Information (SI). Particle dynamics on these potentials were modeled using Langevin and/or Brownian and/or deterministic (Newtonian) dynamics using the internal PyRETIS engine (see SI for parameters). RETIS simulations were performed to serve as reference for the MSM-derived REPPTIS results.

A first set of potentials consists of $M$ cosine-shaped bumps bound by harmonic walls. The height of the bumps is either chosen to be symmetric (Fig. 5A, $M=3$),

11. PAPER III (IN PREPARATION): ESTIMATING FULL PATH
LENGTHS AND KINETICS FROM PARTIAL PATH TRANSITION
INTERFACE SAMPLING SIMULATIONS

9

or asymmetric (Fig. 5B, $M=3$, 'metastable bump'). For the metastable bump potential, interface density and positioning was varied to investigate robustness of the MSM results. This comprises of (a) shifting the interfaces (see SI), or (b) increasing the number of interfaces used (Fig. 5B). Setting $M = 0$, a flat free energy profile is modeled. Setting the barrier height of a cosine bump to a negative number, a cosine dip is obtained. A final potential is a modulated cosine-shaped well, representing a rugged profile with many small metastable states within the well (Fig. 5C).

The results of all the simulations are summarized in Table III, where the estimates for the global crossing probabilities $P_A(\lambda_A \rightarrow \lambda_B)$ and the average $[0^+]$ path lengths $\tau_{[0^+]}$ are given. Both RETIS and MSM results for both properties are seen to be in very well agreement, effectively validating our approach.

## V.   APPLICATION: TRYPSIN-BENZAMIDINE DISSOCATION KINETICS

Static (thermodynamic) properties such as binding affinity (dissociation constant $K_d$) and IC$_{50}$ (the concentration required to cause 50% target inhibition), have long been the primary predictors of drug efficacy [26–29]. The importance of kinetics, in particular the drug residence time, has gained increased recognition for their better correlation with *in vivo* drug efficacy [7, 8, 30]. As drug dissociation is often beyond the timescale accessible to conventional MD, enhanced methods are required, where the trypsin-benzamidine complex has become an exemplar for advanced computational kinetics models [31]. Following this trend, REPPTIS is now applied to compute the trypsin-benzamidine dissociation rate, which is experimentally determined as $k_{\text{off}} = 600\,\text{s}^{-1}$ [32].

REPPTIS requires initial paths for its ensembles before the path sampling simulation can commence. Estimates for these paths are constructed by first running a conventional MD simulation, after which a reactive trajectory is constructed via steered MD. As the initial trajectory is expected not to construct (non-biased) representatives of the path ensembles, an initialization period is required for the importance sampling scheme to generate representative paths. As such, the first $N_{init}$ generated paths are to be excluded from the analysis procedure.

### A.   Preparation for REPPTIS

The MD simulations were performed using Gromacs (version 2021.3) [33]. The Amber14SB force field parameters were used for trypsin and the solvent [34], where the TIP3P water model was used [35]. For benzamidine, the recently developed *ad hoc* parameters of Ref. [36] were used, where the partial charges were obtained using RESP [37]. The non-bonded parameters were obtained
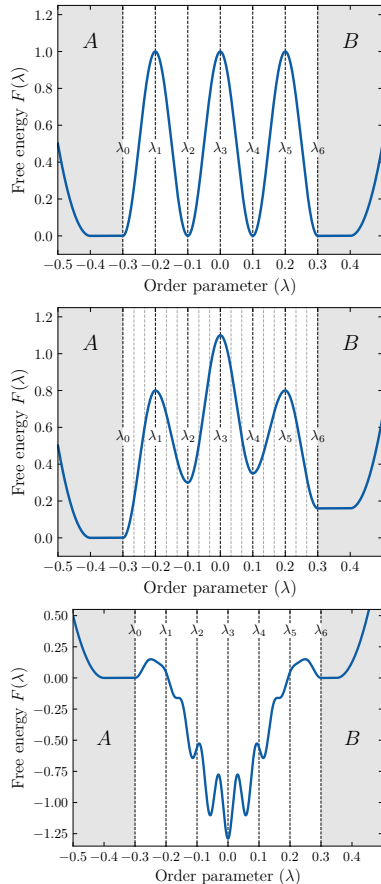


FIG. 5. 1D potential with 3 symmetrical cosine bumps (top), a 'metastable bump' consisting of 3 asymmetric cosine bumps (middle) and a rugged potential (bottom). Stable states $A$ and $B$ and the interfaces are indicated. The middle asymmetric potential (metastable bump) was simulated using the sparse set of interfaces (continuous lines), and a finer set of interfaces (continuous + dashed lines). The analytical expressions for these potentials are given in the SI.

using the Antechamber package [38], after which the parameters of the dihedral bond between the amidine group and benzene were fitted to a QM derived potential energy function. The force field was converted from AMBER to Gromacs format using the ParmED force field conversion tool [39]. The benzamidine molecule has a

| Potential | Dynamics | $P_A(\lambda_B\|\lambda_A)$ | | $\tau_{[0^+]}$ | |
|---|---|---|---|---|---|
| | | RETIS | REPPTIS | RETIS | MSM |
| Flat | Brownian | 0.026 (7%) | 0.026 (8%) | 58.49 | 52.45 |
| | Langevin | 0.516 (2%) | 0.507 (1%) | 2511.8 | 2549.5 |
| | Newtonian | 1.0 (0%) | 1.0 (0%) | 2530.5 | 2428.4 |
| 2 bumps | Brownian | 0.015 (8%) | 0.015 (5%) | 35.0 | 34.36 |
| | Langevin | 0.237 (3%) | 0.227 (2%) | 1583.4 | 1587.9 |
| 3 bumps | Brownian | 0.010 (9%) | 0.010 (6%) | 62.66 | 52.06 |
| 2 dips | Brownian | 0.040 (8%) | 0.040 (5%) | 99.0 | 94.0 |
| Metastable bump | Brownian | 0.011 (9%) | 0.010 (4%) | 58.55 | 49.05 |
| | $\hookrightarrow$ fine | 0.011 (7%) | 0.011 (10%) | 49.44 | 48.50 |
| | $\hookrightarrow$ shifted | 0.009 (6%) | 0.010 (4%) | 51.57 | 50.04 |
| | Langevin | 0.212 (2%) | 0.203 (2%) | 2278.5 | 2243.76 |
| | $\hookrightarrow$ fine | – | 0.170 (7%) | 2278.5 | 2243.76 |
| | $\hookrightarrow$ shifted | – | | | |
| Rugged dip | Brownian | 0.025 (8%) | 0.025 (5%) | 129.710 | 135.118 |

TABLE III. Results of the 1D potential simulations using RETIS and REPPTIS for different dynamics. The times $\tau_{[0^+]}$ are expressed in reduced time units (see SI). The percentages denote the relative errors on the RETIS and REPPTIS crossing probabilities $P_A(\lambda_A \to \lambda_B)$ estimates. An error estimate for the $\tau_{[0^+]}$ MSM results is currently being constructed. The 'fine' and 'shifted' rows of the metastable state denote a finer grid of interfaces and a shifted grid of interfaces, respectively, used for the simulations. Missing entries indicate that the simulations have not (yet) converged.

net charge $q^{\text{net}}$, as the amidine group is protonated at physiological pH ($pk_a = 11.6$, $q_{\text{net}} = +1$) [40]. The system was solvated in a cubic box, ensuring a minimum distance of 1.5 nm between periodic images. Potassium $K^+$ and chlorine $Cl^-$ ions were added to neutralize the system and to reach a physiological salt concentration of 0.15 M. The solvent content consisted of 10590 water molecules, 19 $K^+$ ions, and 28 $Cl^-$ ions. Steepest descent energy minimization was performed until the maximum force was below $1000 \, \text{kj} \, \text{mol}^{-1} \, \text{nm}^{-1}$, after which short equilibration runs in the NVT and NPT ensemble were performed for 100 ps each, where both trypsin and benzamidine heavy atoms were restrained. The unrestrained production simulation ran for 500 ns in the NPT ensemble at 298.15 K using the Nosé-Hoover thermostat [41, 42] (coupling constant of 1 ps) and at 1 bar using the Parinello-Rahman barostat [43] (coupling constant of 5 ps and compressibility of $4.5 \times 10^{-5} \, \text{bar}^{-1}$). A time step of 2 fs was used, where hydrogen bonds were constrained using LINCS [44].

A reactive trajectory was then created using a steered MD simulation, also performed with Gromacs. The center of mass (COM) distance between benzamidine and trypsin was biased with a harmonic potential (force constant $k = 1000 \, \text{kJ} \, \text{mol}^{-1} \, \text{nm}^{-2}$) that was moved at a slow rate of 0.05 nm/ns. The steered MD simulation was performed until the COM distance reached half the box size of $\sim 3.46$ nm in $\sim 40.25$ ns.

The order parameter $\lambda$ for the REPPTIS simulations was then defined as the simple distance metric

$$\lambda(t) = \left\| \mathbf{r}_{C_\gamma^{\text{ASP-189}}}(t) - \mathbf{r}_{C_1^{\text{ben}}}(t) \right\|, \qquad (27)$$

where $\mathbf{r}_{C_\gamma^{\text{ASP-189}}}(t)$ is the position of the $\gamma$ carbon of the aspartic acid (residue 189) and $\mathbf{r}_{C_1^{\text{ben}}}(t)$ is the position of amidine carbon atom of benzamidine (Fig. 6B). The REPPTIS simulation consisted of 33 ensembles $\{E_i\}_{i=0}^{32} = \{[0^-], [0^{\pm\prime}], [1^\pm], \ldots, [31^\pm]\}$, where 33 interfaces were positioned at $\{\lambda_i\}_{i=0}^{32} = [4, 4.15, 4.3, 4.45, 4.6, 4.8, 5, 5.33, 5.66, 6, 6.33, 6.66, 7, 7.25, 7.5, 7.75, 8, 8.33, 8.66, 9, 9.33, 9.66, 10, 10.5, 11, 12, 13, 14, 15, 16, 17, 18, 19]$ Å. Interface positioning is most dense for small $\lambda$ values where the free energy profile is expected to rise sharply. Order parameters were calculated every 20 fs, and trial MC moves were set at 50 % shooting moves and 50 % replica exchange moves. The simulation was performed with a customized hybrid version of $\infty$RETIS [45, 46] and PyRETIS 3 [47]. This was done for the ability to use more hardware (asynchronous formalism), where the infinite swapping formalism does not apply to REPPTIS.

## B. Results

The REPPTIS simulation produced 628 682 MC moves for a total of 1.079 μs MD simulation time. Of this simulation time, 0.878 μs was performed by shooting moves (of which 0.568 μs was accepted) and 0.201 μs by replica exchange moves (automatically accepted when performed). The first $N_{\text{init}} = 100000$ paths were discarded for the analysis to avoid initialization effects. An exception to this is the running estimate of the global crossing probability $P_A(\lambda_A|\lambda_B)$ in Fig. 6D.
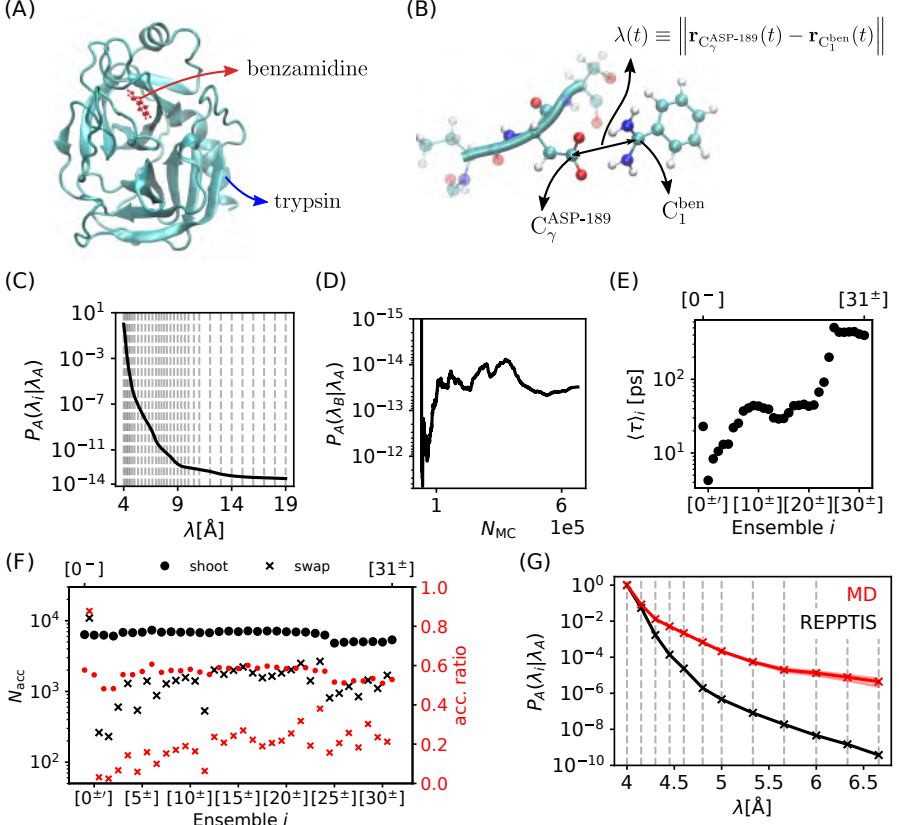
FIG. 6. **A**: The trypsin-benzamidine complex. **B**: the order parameter $\lambda$ is the distance between the $\gamma$ carbon atom of ASP-189 and the amidine carbon of benzamidine. **C**: The crossing probability profile $P_A(\lambda_i|\lambda_A)$. **D**: The running estimate of the global crossing probability $P_A(\lambda_B|\lambda_A)$. **E**: The average path lengths $\langle\tau\rangle_i$ of the ensembles $E_i = \{[0^-], [0^{\pm'}], [1^\pm], \ldots, [31^\pm]\}$. **F**: Statistics of the REPPTIS simulation. Black dots represent the amount of accepted shooting moves in the ensembles, while black crosses represent the amount of accepted swapping moves between neighboring ensembles. In red, the acceptance ratio of shooting (dots) and swap (crosses) trials is shown. **G**: The $P_A(\lambda_i|\lambda_A)$ profile in the $\lambda \in [4, 6.66]$ Å range, as estimated by a 450 ns MD trajectory (red) and REPPTIS (black). Red shade indicates the (Poisson) std. err. of first crossing counts $N_{\text{counts}}$ in the MD trajectory (rel. err. of $1/\sqrt{N_{\text{counts}}}$).

The final estimate for the global crossing probability was $P_A(\lambda_B|\lambda_A) = 2.47 \times 10^{-14} \pm 79\%$, with a positive flux $f_A = 2.17 \text{ps}^{-1}$, resulting in a dissociation rate $k_{\text{off}} = 0.05\,\text{s}^{-1} \pm 79\%$. This result is a significant underestimation of the experimental value and other computational works [31]. It was assumed that the ensembles close to the bound state were not sampled efficiently, which is now discussed.

The average path lengths of the first 10 positive en-

sembles are small (Fig. 6E), which is expected due to attractive forces at the binding pocket. If large energy barriers orthogonal to the $\lambda$ parameter are present, the shooting moves will likely not overcome them, resulting in localized sampling of path space. REPPTIS then depends on the replica exchange formalism to allow these ensembles to explore (more) favorable regions. The acceptance ratio for the shooting moves and replica exchange moves are shown in Fig. 6F, where the acceptance

ratio of $[0^{\pm\prime}] \leftrightarrow [1^{\pm}]$ (3.1 %) and $[1^{\pm}] \leftrightarrow [2^{\pm}]$ (2.5 %) exchange are seen to be especially low. Low swap acceptance ratios are expected for ensembles distributed over a steep free energy profile, where most paths are of type LML , providing no path overlap between adjacent ensembles required for replica exchange. This is problematic if the initial paths of these ensembles are not representative of the underlying path ensembles, where (a) the shooting move has trouble overcoming orthogonal barriers, and (b) the replica exchange move has trouble finding path overlap. As the steered MD bias was performed on the distance between the trypsin COM and the benzamidine COM rather than the distance of benzamidine to the center of the binding pocket, the initial path contained a directional bias, which may have pulled benzamidine along an unfavorable dissociation pathway. While ensembles far from the bound state could relax to more favorable regions, the ensembles close to the bound state could not. This was further complicated by the occurrence of LINCS warnings for benzamidine H-bonds in the $[2^{\pm}]$ paths, which cut the simulation short. Lowering the simulation time step to 1 fs did not resolve LINCS errors, and it is not clear whether the *ad hoc* force field parameters or an unfortunate chain of shooting moves caused this issue.

As these problematic ensembles are close to the bound state, it was possible to calculate the crossing probability profile $P_A(\lambda_i|\lambda_A)$ for the $\lambda \in [4, 6.66]$ Å range using a 450 ns brute-force MD simulation where the $\lambda$ parameter was saved every 20 fs. The difference between the REPPTIS and MD profiles is shown in Fig. 6G. REPPTIS underestimates the $P_A(\lambda_i|\lambda_A)$ profile by approximately 4 orders of magnitude. Appending the $\lambda > 6.66$ Å crossing probability profile of REPPTIS to the $\lambda < 6.66$ Å profile of the MD simulation results in a rate constant $k_{\text{off}}^{\text{MD+REPPTIS}} = 635\,\text{s}^{-1}$ that lies very close to the experimental value.

The largest discrepancy between MD and REPPTIS occurs in the $[4.15 < \lambda < 4.8]$ Å region, corresponding to ensembles $\{[1^{\pm}], \ldots, [4^{\pm}]\}$. Crossing of the $\lambda_1 = 4.15$ Å interface is, however, accurately predicted by REPPTIS, as the $[0^-] \leftrightarrow [0^{\pm\prime}]$ exchange is highly accepted. The long MD simulation allowed direct estimation of the flux, where $f_A^{\text{MD}} = 2.01\,\text{ps}^{-1}$ is in good agreement with the MSM-derived value of $f_A = 2.17\,\text{ps}^{-1}$.

There is clearly a need for the REPPTIS methodology to better sample metastable states separated by barriers orthogonal to $\lambda$. While a combination of better path initialization and better choice of $\lambda$ parameter may have significantly improved sampling close to the bound state, this need remains due to biological systems often being more complex than the trypsin-benzamidine system considered here. An extension of REPPTIS to a multi-dimensional order parameter space is therefore desirable, where an enhanced free energy sampling proce-

dure can be used to first determine relevant collective variables along the reactive pathway(s). Another possibility to tackle specifically steep energy regions is to construct a hybrid RETIS and REPPTIS methodology. RETIS-like ensembles could be positioned in the steep energy regions and connected to REPPTIS ensembles elsewhere. A methodology to exchange paths between these different ensembles should then be developed.

## VI. CONCLUSION

In conclusion, we have introduced a MSM analysis framework that is compatible with the memory assumptions of (RE)PPTIS path ensembles. This new approach enables the calculation of time-dependent properties such as MFPTs, fluxes, and rates, effectively addressing a significant limitation of REPPTIS. We validated our framework using one-dimensional potential systems with various dynamics, demonstrating consistency with RETIS results.

Application of REPPTIS to the trypsin-benzamidine system revealed that orthogonal barriers separating metastable states challenges efficient sampling. By combining REPPTIS with brute-force MD simulations, a dissociation rate that aligns closely with experimental data could still be recovered. These findings highlight the need for further improvements, such as exploring multidimensional REPPTIS path ensemble definitions or hybrid methods that combine TIS and PPTIS ensembles. Such enhancements would facilitate more efficient sampling of biological systems, which often exhibit greater complexity than the trypsin-benzamidine system.

Overall, our MSM-based analysis framework significantly extends the capabilities of REPPTIS, providing a robust tool for investigating the kinetics of rare and slow events in molecular systems, and opening new avenues for studying biomolecular mechanisms and drug kinetics.

## SUPPORTING INFORMATION

The Supporting Information is available . . .

## CODE AVAILABILITY

The MSM analysis framework is available on GitHub at https://github.com/annekegh/tistools.

## ACKNOWLEDGMENTS

# 11. Paper III (in preparation): Estimating full path lengths and kinetics from partial path transition interface sampling simulations

13

[1] R. O. Dror, R. M. Dirks, J. Grossman, H. Xu, and D. E. Shaw, Annual review of biophysics **41**, 429 (2012).

[2] S. A. Hollingsworth and R. O. Dror, Neuron **99**, 1129 (2018).

[3] D. E. Shaw, P. J. Adams, A. Azaria, J. A. Bank, B. Batson, A. Bell, M. Bergdorf, J. Bhatt, J. A. Butts, T. Correia, R. M. Dirks, R. O. Dror, M. P. Eastwood, B. Edwards, A. Even, P. Feldmann, M. Fenn, C. H. Fenton, A. Forte, J. Gagliardo, G. Gill, M. Gorlatova, B. Greskamp, J. Grossman, J. Gullingsrud, A. Harper, W. Hasenplaugh, M. Heily, B. C. Heshmat, J. Hunt, D. J. Ierardi, L. Iserovich, B. L. Jackson, N. P. Johnson, M. M. Kirk, J. L. Klepeis, J. S. Kuskin, K. M. Mackenzie, R. J. Mader, R. McGowen, A. McLaughlin, M. A. Moraes, M. H. Nasr, L. J. Nociolo, L. O'Donnell, A. Parker, J. L. Peticolas, G. Pocina, C. Predescu, T. Quan, J. K. Salmon, C. Schwink, K. S. Shim, N. Siddique, J. Spengler, T. Szalay, R. Tabladillo, R. Tartler, A. G. Taube, M. Theobald, B. Towles, W. Vick, S. C. Wang, M. Wazlowski, M. J. Weingarten, J. M. Williams, and K. A. Yuh, in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, SC '21 (Association for Computing Machinery, New York, NY, USA, 2021).

[4] R. C. Bernardi, M. C. Melo, and K. Schulten, Biochimica et Biophysica Acta (BBA)-General Subjects **1850**, 872 (2015).

[5] M. C. Zwier and L. T. Chong, Current opinion in pharmacology **10**, 745 (2010).

[6] K. Henzler-Wildman and D. Kern, Nature **450**, 964 (2007).

[7] R. A. Copeland, D. L. Pompliano, and T. D. Meek, Nature reviews Drug discovery **5**, 730 (2006).

[8] R. A. Copeland, Nature Reviews Drug Discovery **15**, 87 (2016).

[9] M. De Vivo, M. Masetti, G. Bottegoni, and A. Cavalli, Journal of medicinal chemistry **59**, 4035 (2016).

[10] D. A. Schuetz, W. E. A. de Witte, Y. C. Wong, B. Knasmueller, L. Richter, D. B. Kokh, S. K. Sadiq, R. Bosma, I. Nederpelt, L. H. Heitman, *et al.*, Drug Discovery Today **22**, 896 (2017).

[11] T. S. Van Erp, D. Moroni, and P. G. Bolhuis, The Journal of chemical physics **118**, 7762 (2003).

[12] T. S. van Erp, Europhysics Letters **143**, 30001 (2023).

[13] T. S. van Erp, The Journal of chemical physics **125** (2006).

[14] T. S. Van Erp, Advances in Chemical Physics **151**, 27 (2012).

[15] R. Cabriolu, K. M. Skjelbred Refsnes, P. G. Bolhuis, and T. S. van Erp, The Journal of Chemical Physics **147** (2017).

[16] D. Moroni, P. G. Bolhuis, and T. S. van Erp, The Journal of chemical physics **120**, 4055 (2004).

[17] W. Vervust, D. T. Zhang, T. S. Van Erp, and A. Ghysels, Biophysical Journal **122**, 2960 (2023).

[18] P. Májek and R. Elber, Journal of chemical theory and computation **6**, 1805 (2010).

[19] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller, J. Chem. Phys. **21**, 1087 (1953).

[20] E. Riccardi, O. Dahlen, and T. S. van Erp, The Journal of Physical Chemistry Letters **8**, 4456 (2017).

[21] A. Ghysels, S. Roet, S. Davoudi, and T. S. van Erp, Physical Review Research **3**, 033068 (2021).

[22] D. T. Zhang, E. Riccardi, and T. S. van Erp, The Journal of Chemical Physics **158** (2023).

[23] R. J. Allen, C. Valeriani, and P. R. Ten Wolde, Journal of physics: Condensed matter **21**, 463102 (2009).

[24] T. S. Van Erp and P. G. Bolhuis, Journal of computational Physics **205**, 157 (2005).

[25] D. Moroni, T. S. van Erp, and P. G. Bolhuis, Physical Review E **71**, 056709 (2005).

[26] C. A. Lipinski, F. Lombardo, B. W. Dominy, and P. J. Feeney, Advanced drug delivery reviews **64**, 4 (2012).

[27] W. L. Jorgensen, Science **303**, 1813 (2004).

[28] R. Claveria-Gimeno, S. Vega, O. Abian, and A. Velazquez-Campoy, Expert Opinion on Drug Discovery **12**, 363 (2017).

[29] D. L. Mobley and M. K. Gilson, Annual review of biophysics **46**, 531 (2017).

[30] P. J. Tummino and R. A. Copeland, Biochemistry **47**, 5481 (2008).

[31] F. Sohraby and A. Nunes-Alves, Trends in Biochemical Sciences **48**, 437 (2023).

[32] F. Guillain and D. Thusius, Journal of the American Chemical Society **92**, 5534 (1970).

[33] M. J. Abraham, T. Murtola, R. Schulz, S. Páll, J. C. Smith, B. Hess, and E. Lindahl, SoftwareX **1**, 19 (2015).

[34] J. A. Maier, C. Martinez, K. Kasavajhala, L. Wickstrom, K. E. Hauser, and C. Simmerling, Journal of chemical theory and computation **11**, 3696 (2015).

[35] W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey, and M. L. Klein, The Journal of chemical physics **79**, 926 (1983).

[36] S. Raniolo and V. Limongelli, Frontiers in molecular biosciences **8** (2021).

[37] C. I. Bayly, P. Cieplak, W. Cornell, and P. A. Kollman, The Journal of Physical Chemistry **97**, 10269 (1993).

[38] J. Wang, R. M. Wolf, J. W. Caldwell, P. A. Kollman, and D. A. Case, Journal of computational chemistry **25**, 1157 (2004).

[39] M. R. Shirts, C. Klein, J. M. Swails, J. Yin, M. K. Gilson, D. L. Mobley, D. A. Case, and E. D. Zhong, Journal of computer-aided molecular design **31**, 147 (2017).

[40] P. Y. Lam, C. G. Clark, R. Li, D. J. Pinto, M. J. Orwat, R. A. Galemmo, J. M. Fevig, C. A. Teleha, R. S. Alexander, A. M. Smallwood, *et al.*, Journal of medicinal chemistry **46**, 4405 (2003).

[41] S. Nosé, Molecular physics **52**, 255 (1984).

[42] W. G. Hoover, Physical review A **31**, 1695 (1985).

[43] M. Parrinello and A. Rahman, Journal of Applied physics **52**, 7182 (1981).

[44] B. Hess, H. Bekker, H. J. Berendsen, and J. G. Fraaije, Journal of computational chemistry **18**, 1463 (1997).

[45] S. Roet, D. T. Zhang, and T. S. van Erp, The Journal of Physical Chemistry A **126**, 8878 (2022).

[46] D. T. Zhang, L. Baldauf, S. Roet, A. Lervik, and T. S. van Erp, Proceedings of the National Academy of Sciences **121**, e2318731121 (2024).

[47] W. Vervust, D. T. Zhang, A. Ghysels, S. Roet, T. S. van Erp, and E. Riccardi, Journal of Computational Chem-

istry (2024).

Supporting information:

Estimating full path lengths and kinetics from partial path transition interface sampling simulations

Wouter Vervust, Elias Wils, and An Ghysels

*IBiTech - BioMMedA group, Ghent University,*

*The Core, Corneel Heymanslaan 10, 9000 Gent, Belgium*

1

**CONTENTS**

## I. 1D POTENTIALS

## II. EQUATIONS FOR PROBABILITIES IN MSM NETWORK

Consider a Markov state model with finite state space $\mathcal{S} = \{S_\alpha, S_\beta, \ldots\}$ and transition matrix $M$, whose elements $M_{\alpha\beta}$ denote the probability to transition from state $S_\alpha$ to state $S_\beta$.

# 11. Paper III (in preparation): Estimating full path lengths and kinetics from partial path transition interface sampling simulations

### A. Probability to hit state $S_\beta$

Let us draft equations for the probability $Q$ to hit $S_\beta$ when your current state is $S_\delta$,

$$Q_{\delta(\beta)} \equiv P(\exists n \geq 0 : S_n = \beta | S_0 = \delta) \tag{1}$$

which is a so-called hitting probability. $\beta$ is fixed in this paragraph. When starting from state $\delta = \beta$, you reach $\beta$ with probability is 1, because you are already in $\beta$,

$$Q_{\beta(\beta)} = 1 \tag{2}$$

When starting from $\delta \neq \beta$, you are not yet in $\beta$, and you will have to take at least one step. By conditioning on this first step, the equations for the other hitting probabilities $Q_{\delta(\beta)}$ values may be constructed,

$$Q_{\delta(\beta)} = P(\exists n \geq 0 : S_n = \beta | S_0 = \delta) \tag{3}$$

$$= P(\exists n \geq 1 : S_n = \beta | S_0 = \delta) \tag{4}$$

$$= \sum_\gamma M_{\delta\gamma} P(\exists n \geq 1 : S_n = \beta | S_0 = \delta, S_1 = \gamma) \tag{5}$$

$$= \sum_\gamma M_{\delta\gamma} P(\exists n \geq 1 : S_n = \beta | S_1 = \gamma) \tag{6}$$

$$= \sum_\gamma M_{\delta\gamma} P(\exists n \geq 0 : S_n = \beta | S_0 = \gamma) \tag{7}$$

$$= \sum_\gamma M_{\delta\gamma} Q_{\gamma(\beta)} \tag{8}$$

The first equality is the definition. The second means that we need to make at least one step because $\delta \neq \beta$. In the third equality, we have conditioned the probability on the outcome of the first step: a transition to any state $S_\gamma$. In the fourth, the Markovian property removes the dependence on the initial state $S_0 = S_\delta$. In the fifth, we can see that starting from $S_1$ is equivalent to start over a new chain so we reset the counter of the steps to 0, and in the last equality we recognize the definition of $Q_{\gamma(\beta)}$.

To solve the set of Eqs. 2-8, an adapted transition matrix $M'$ is constructed where state $\beta$ is absorbing by adapting row $\beta$:

$$M'_{\beta\beta} = 1; \; M'_{\beta,\gamma} = 0, \forall \gamma \neq \beta \tag{9}$$

3

These matrix elements do not contribute to Eq. 8, so adapting the elements lets us write a compact version of Eq. 8. Indeed, given a specific $\beta$, the equations for the $Q_{(\beta)}$ vector can be summarized, in combination with Eq. 2, by

$$Q_{(\beta)} = M'Q_{(\beta)} \tag{10}$$

See section IV on how to solve this matrix equation.

### B. Probability to hit state $S_\beta$ before state $S_\alpha$

In a next step, we compute the probability $Z$ that a trajectory starting in state $S_\delta$ can reach state $S_\beta$ before reaching state $S_\alpha$,

$$Z_{\delta(\beta)} = P(\exists n \geq 0 : S_n = \beta \,\text{before}\, S_n = \alpha | S_0 = \delta) \tag{11}$$

which is also a hitting probability. We draft the equations for different starting values $\delta$, where the states $\alpha, \beta$ are fixed in this paragraph.

Starting from state $\delta = \beta$, you are already in the desired state with certainty, so the probability is 1. However, starting from $\delta = \alpha$, it is certain that you will not reach $\beta$, because you already reached the $\alpha$ boundary. This gives the following $Z$ values for these boundaries,

$$Z_{\beta(\beta)} = 1 \tag{12}$$

$$Z_{\alpha(\beta)} = 0 \tag{13}$$

When starting from $\delta \neq \beta$ and $\delta \neq \alpha$, you need to take at least one step ($n \geq 1$). By conditioning the probability on this first step, the equations become

$$Z_{\delta(\beta)} = P(\exists n \geq 0 : S_n = \beta \,\text{before}\, S_n = \alpha | S_0 = \delta) \tag{14}$$

$$= P(\exists n \geq 1 : S_n = \beta \,\text{before}\, S_n = \alpha | S_0 = \delta) \tag{15}$$

$$= \sum_\gamma M_{\delta\gamma} P(\exists n \geq 1 : S_n = \beta \,\text{before}\, S_n = \alpha | S_0 = \delta, S_1 = \gamma) \tag{16}$$

$$= \sum_\gamma M_{\delta\gamma} P(\exists n \geq 1 : S_n = \beta \,\text{before}\, S_n = \alpha | S_1 = \gamma) \tag{17}$$

$$= \sum_\gamma M_{\delta\gamma} P(\exists n \geq 0 : S_n = \beta \,\text{before}\, S_n = \alpha | S_0 = \gamma) \tag{18}$$

$$= \sum_\gamma M_{\delta\gamma} Z_{\gamma(\beta)} \tag{19}$$

4

where a similar reasoning is followed as in the derivation of Eq. 8. These equations can be
implemented in practice with an adapted transition matrix $M''$ where the two states $\alpha$ and
$\beta$ are both absorbing by adapting their rows,

$$M''_{\alpha\alpha} = 1; M''_{\alpha,\gamma} = 0, \forall \gamma \neq \alpha \tag{20}$$

$$M''_{\beta\beta} = 1; M''_{\beta,\gamma} = 0, \forall \gamma \neq \beta \tag{21}$$

The equations for the $Z_{(\beta)}$ vector with probabilities to reach the product state from the
several states, before reaching the reactant state, can then be summarized as

$$Z_{(\beta)} = M'' Z_{(\beta)} \tag{22}$$

in combination with Eqs. 12-13. See again section IV on how to solve this matrix equation.

### C. Probability to hit state $S_\beta$ before state $S_\alpha$, after leaving state $S_\alpha$

Next, we compute the probability $Y$ that a trajectory that left state $S_\alpha$, can reach state
$S_\beta$ before returning to the initial state $S_\alpha$,

$$Y_{\alpha(\beta)} = P(\exists n \geq 1 : S_n = \beta \, \text{before} \, S_n = \alpha | S_0 = \alpha) \tag{23}$$

This probability is the complement of the *return* probability to state $\alpha$. The number of
steps must be at least $n = 1$. It is assumed that $\alpha \neq \beta$.

By conditioning on the first step, the probability to reach $\beta$ before returning to $\alpha$ can be
written in terms of the previous $Z$ vectors,

$$Y_{\alpha(\beta)} = P(\exists n \geq 1 : S_n = \beta \, \text{before} \, S_n = \alpha | S_0 = \alpha) \tag{24}$$

$$= \sum_\gamma M_{\alpha\gamma} P(\exists n \geq 1 : S_n = \beta \, \text{before} \, S_n = \alpha | S_0 = \alpha, S_1 = \gamma) \tag{25}$$

$$= \sum_\gamma M_{\alpha\gamma} P(\exists n \geq 1 : S_n = \beta \, \text{before} \, S_n = \alpha | S_1 = \gamma) \tag{26}$$

$$= \sum_\gamma M_{\alpha\gamma} P(\exists n \geq 0 : S_n = \beta \, \text{before} \, S_n = \alpha | S_0 = \gamma) \tag{27}$$

$$= \sum_\gamma M_{\alpha\gamma} Z_{\gamma(\beta)} \tag{28}$$

with similar justifications for the equalities as for the $Q$ or $Z$ vectors in Eqs. 8 or 19,
respectively.

5

If we let the starting state vary, this reads in matrix notation as

$$Y_{(\beta)} = M Z_{(\beta)} \tag{29}$$

where the computation of the $Y_{\beta(\beta)}$ component does not need to be executed as this has no meaning. By definition, $Y_{\delta(\beta)}$ for $\delta \notin \{\alpha, \beta\}$ is equal to $Z_{\delta(\beta)}$.

## III.   MSM EQUATIONS FOR TIMES

### A.   Average time to hit state $S_\beta$

We first quickly recap what is meant by accumulated time. States in the Markov model represent path segments, and chains represent (overlapping) longer paths that are built by concatenating these segments. As such, a chain of $n$ states represents a path of $n$ segments, where the accumulated time represents the total length of the $n$ segments. The overlap between segments was accounted for by only considering the $\tau_{(m2)}$ parts of the segment path lengths $\tau = \tau_{(1m2)}$ (see main text). The accumulated path of a chain thus relates to the sum of these $\tau_{(m2)}$ parts of its corresponding segments. Depending on the property of interest, the starting and/or ending segments may incorporate or lose their $\tau_{(1)}$ or $\tau_{(m2)}$ part.

Assume $F_{\delta(\beta)}$ is the average accumulated time to reach state $\beta$ starting from state $\delta$. When starting in $\delta = \beta$ itself, the sequence $(U_0, \ldots)$ is stopped immediately and it consists of only one state $(\beta)$. The accumulated time for this boundary is then the (m2) parts of state $\beta$ itself,

$$F_{\beta(\beta)} = \tau_{(m2),\beta} \tag{30}$$

For $\delta \neq \beta$, at least one step needs to be taken. By conditioning on this step, a recursive relation can be built (with $E$ denoting the expected value),

$$F_{\delta(\beta)} = E(t|S_0 = \delta) \tag{31}$$

$$= \sum_\gamma M_{\delta\gamma} E(t|S_0 = \delta, S_1 = \gamma) \tag{32}$$

$$= \sum_\gamma M_{\delta\gamma} (\tau_{(m2),\delta} + E(t|S_0 = \gamma)) \tag{33}$$

$$= \tau_{(m2),\delta} + \sum_\gamma M_{\delta\gamma} F_{\gamma(\beta)} \tag{34}$$

6

## 11. PAPER III (IN PREPARATION): ESTIMATING FULL PATH LENGTHS AND KINETICS FROM PARTIAL PATH TRANSITION INTERFACE SAMPLING SIMULATIONS

The time length of a single step is again measured in accumulated time, which gives the $\tau_{(m2),\delta}$ term. In the last equality, the probability conservation $\sum_\gamma M_{\delta\gamma} = 1$ was used, and the definition of $F_{\gamma(\beta)}$ was recognized. Together with Eq. 30, these equations can be written in matrix notation,

$$F_{(\beta)} = \tau'_{(m2)} + M'_\dagger F_{(\beta)} \tag{35}$$

Here, $M'_\dagger$ is the similar to the adapted transition matrix $M'$ as in Eq. 10, but where now both the row and column corresponding to state $\beta$ are made adsorbing. In addition, $\tau'_{(m2)}$ is a column vector that contains the $\tau_{(m2),\delta}$ elements with however a zero entry for the $\beta$ element of this vector.

### B. Average time to hit state $S_\alpha$ or $S_\beta$

Next, assume $G_{\delta(\beta)}$ is the average accumulated time to reach state $\alpha$ or $\beta$ starting from state $\delta$. Assume the reactant state $A$ is the set with adsorbing states (here, $A = \{\alpha\}$), the product state $B$ is the second set with adsorbing states (here, $B = \{\beta\}$), and the set $C = A \cup B$.

$$G_{\delta(C)} \equiv E(t|S_0 = \delta) \tag{36}$$

For $\delta \in C$, the boundaries are

$$G_{\delta(C)} = \tau_{(m2),\delta} \tag{37}$$

For $\delta \notin C$, at least one step needs to be taken. By conditioning on this first step, a recursive relation can be built,

$$G_{\delta(C)} = E(t|S_0 = \delta) \tag{38}$$

$$= \sum_\gamma M_{\delta\gamma} E(t|S_0 = \delta, S_1 = \gamma) \tag{39}$$

$$= \sum_\gamma M_{\delta\gamma} (\tau_{(m2),\delta} + E(t|S_0 = \gamma)) \tag{40}$$

$$= \tau_{(m2),\delta} + \sum_\gamma M_{\delta\gamma} G_{\gamma(C)} \tag{41}$$

The time length of a single step is again measured in accumulated time, which gives the $\tau_{(m2),\delta}$ term. Together with Eq. 37, this can be written in matrix notation

$$G_{(C)} = \tau''_{(m2)} + M'' G_{(C)} \tag{42}$$

7

158

where $\tau''$ is the column vector with the $\tau_{(m2),\delta}$ elements and zero entries for elements corresponding to $C$ states, and $M''$ is an adapted transition matrix $M$ where all rows corresponding to states in $C$ have been made adsorbing.

## C. Average time to hit state $S_\alpha$ or $S_\beta$, after leaving state $S_\alpha$

Next, assume $H_{\delta(\beta)}$ is the average accumulated time to reach state $\alpha$ or $\beta$ starting from state $\alpha$, given that at least one step has been taken to leave $\alpha$. Conditioning on this first step gives an expected stopping time $H_{\alpha(C)}$,

$$H_{\alpha(C)} \equiv E(t, n \geq 1 | S_0 = \alpha) \tag{43}$$

$$= \sum_\gamma M_{\alpha\gamma} E(t, n \geq 1 | S_0 = \alpha, S_1 = \gamma) \tag{44}$$

$$= \sum_\gamma M_{\alpha\gamma} (\tau_{(m2),\alpha} + E(t, n \geq 0 | S_0 = \gamma)) \tag{45}$$

$$= \tau_{(m2),\alpha} + \sum_\gamma M_{\alpha\gamma} G_{\gamma(C)} \tag{46}$$

$$= \tau_{(m2),\alpha} + \left( M G_{(C)} \right)_\alpha \tag{47}$$

where we wrote the matrix product $M G_{(C)}$. Using the previously computed $G_{(\beta)}$ vector, $H_{\alpha(C)}$ can be computed as well with Eq. 47.

## IV. HOW TO SOLVE THE EQUATIONS

### A. Solving equations for probabilities

In order to solve Eq. 10 (Eq. 22) for the vector $Q_{(\beta)}$ (vector $Z_{(\beta)}$), the matrix $M$ is reorganized to put the absorbing rows for $\beta$ (for $\alpha$ and $\beta$) to the top and corresponding columns to the left of the matrix. For generality, assume the set $C$ contains all $n_C$ absorbing states, and $n_s$ is the total number of states. The reorganized matrix has four blocks,

$$M = \begin{pmatrix} M_C & E \\ D & \tilde{M} \end{pmatrix} \tag{48}$$

The diagonal block $M_C$ has dimension $n_C \times n_C$, and the diagonal block $\tilde{M}$ has dimension $(n_s - n_C) \times (n_s - n_C)$. The off-diagonal block $D$ has dimension $(n_s - n_C) \times n_C$, while $E$ has dimension $n_C \times (n_s - n_C)$.

8

## 11. PAPER III (IN PREPARATION): ESTIMATING FULL PATH LENGTHS AND KINETICS FROM PARTIAL PATH TRANSITION INTERFACE SAMPLING SIMULATIONS

The equations for $Q$ or $Z$ are of the shape $q = M'q$, where $M'$ is the adapted transition matrix where the $M_C$ and $E$ block in the transition matrix $M$ are set to adsorbing rows. Introducing $q_C$ as the part of $q$ that corresponds to the absorbing states in $C$ and $\tilde{q}$ for the remainder of the $q$ vector, the equations $q = M'q$ read in block diagonal matrix form,

$$\begin{pmatrix} q_C \\ \tilde{q} \end{pmatrix} = \begin{pmatrix} 1_{n_C} & 0 \\ D & \tilde{M} \end{pmatrix} \begin{pmatrix} q_C^0 \\ \tilde{q} \end{pmatrix} \tag{49}$$

with $1_{n_C}$ an identity matrix of dimension $n_C$ and $q_C^0$ are the known boundary values. The matrix equation is equivalent to

$$\begin{cases} q_C &= q_C^0 \\ \tilde{q} &= Dq_C^0 + \tilde{M}\tilde{q} \end{cases} \tag{50}$$

with solution

$$\begin{cases} q_C &= q_C^0 \\ \tilde{q} &= (1 - \tilde{M})^{-1} Dq_C^0 \end{cases} \tag{51}$$

Next, a quantity of the shape $y = Mq$ can be computed, for instance for the computation of $Y$ in Eq. 28. Using a similar decomposition of $y$ into $y_C$ and $\tilde{y}$, the vector $y$ becomes

$$\begin{pmatrix} y_C \\ \tilde{y} \end{pmatrix} = \begin{pmatrix} M_C & E \\ D & \tilde{M} \end{pmatrix} \begin{pmatrix} q_C^0 \\ \tilde{q} \end{pmatrix} \tag{52}$$

and thus

$$\begin{cases} y_C &= M_C q_C^0 + E(1 - \tilde{M})^{-1} Dq_C^0 \\ \tilde{y} &= Dq_C^0 + \tilde{M}(1 - \tilde{M})^{-1} Dq_C^0 \end{cases} \tag{53}$$

where the interesting elements are located in $y_C$.

### B. Solving MFPT equations

To solve for $F_{(\beta)}$ in Eq. 35, the matrix $M'_\dagger$ is reorganized similarly as before, where now the column corresponding to state $\beta$ is absorbing,

$$F_{(\beta)} = \begin{pmatrix} F_{\beta(\beta)} \\ \tilde{F}_{(\beta)} \end{pmatrix} = \begin{pmatrix} 0 \\ \tilde{\tau}_{(m2)} \end{pmatrix} + \begin{pmatrix} 1 & 0 \\ 0 & \tilde{M} \end{pmatrix} \begin{pmatrix} \tau_{(m2),\beta} \\ \tilde{F}_{(\beta)} \end{pmatrix}, \tag{54}$$

with $\tau'_{(m2)} = (0, \tilde{\tau}_{(m2)})^T$ and $F_{(\beta)} = (\tau_{(m2),\beta}, \tilde{F}_{(\beta)})^T$. The solution of the interesting elements $\tilde{F}_{(\beta)}$ is

$$F_{(\beta)} = (1 - \tilde{M})^{-1} \tilde{\tau}_{(m2)}. \tag{55}$$

9

160

For $G_{(C)}$ in Eq. 42, the matrix $M''_\dagger$ is reorganized similarly, where the absorbing columns/rows now have dimension $n_C$. Denoting $G_{(C)}$ as $(\tau_{(m2),\delta\in C}, \tilde{G}_{(C)})$ and $\tau''_{(m2)}$ as $(0, \tilde{\tau}_{(m2)})$, the matrix equation reads

$$G_{(C)} = \begin{pmatrix} 0 \\ \tilde{\tau}_{(m2)} \end{pmatrix} + \begin{pmatrix} 1_{n_C} & 0 \\ 0 & \tilde{M} \end{pmatrix} \begin{pmatrix} \tau_{(m2),\delta\in C} \\ \tilde{G}_{(C)} \end{pmatrix} \tag{56}$$

with solution for $\tilde{G}_{(C)}$ as

$$\tilde{G}_{(C)} = \left(1 - \tilde{M}\right)^{-1} \tilde{\tau}_{(m2)} \tag{57}$$

For $H_{\alpha(C)}$ in Eq. 47, the matrix equation was already written in a form that can be solved directly.


## V. ONE-DIMENSIONAL POTENTIALS

### A. Simple, symmetric 1D potentials with intuitive interface placement

We have constructed several potentials for which we simulate a single one-dimensional Langevin particle, which are based on Ref. [1]. The potentials $V(x)$ consist of one or more cosine-shaped bumps bound by two harmonic walls on both sides, centered around $x = 0$. The reactant state $A$ and product state $B$ are located on either side of the bump(s). These are described by a piecewise function of the particle's position $x$:

$$V(x) = \begin{cases} \frac{1}{2}V_0 \left( \cos \frac{\pi(x-a(M+1))}{a} + 1 \right), & |x| \leq Ma \\ 0, & Ma < |x| \leq b \\ \frac{1}{2}k(|x| - b)^2. & |x| \geq b \end{cases} \tag{58}$$

The force on the particle derived from this potential is continuous everywhere. The parameter $V_0$ represents the bump(s)'s height, $k$ is the strength of the harmonic walls, and $M$ is the number of cosine-shaped bumps with period $2a$. Setting $V_0 = 0$ results in a flat potential function. For our application, $a = 0.1$, $b = Ma + 0.1$, $k = 100$ and $V_0 = 1$ (or 0 for a flat potential).

The integration of the equations of motion was done by the internal PyRETIS 3 engine [2], where reduced units are used in a Lennard-Jones type unit system based on argon [3]. The order parameter $\lambda$ is the position $x$.

10

For the simple cosine bumps (Eq. (58)) the parameters were set according to Ref. [1]
with a friction coefficient $\gamma = 5$ when using Langevin dynamics. The number of cycles was
chosen to be $30\,000$, and the number of interfaces $N$ was chosen proportional to the number
of bumps $(2M+1)$. They were uniformly spaced with a distance $a$ between them, such that
they correspond with the extrema of the cosine bumps. With the flat potential $(V_0 = 0)$ we
imposed 5 equally-spaced interfaces.

### B.   More complex 1D potentials with varied interface placement

Two more potentials containing metastable states were constructed. On one hand, we
constructed an asymmetric 1D potential containing three smaller cosine-shaped bumps on
top of a bigger one (metastable bump) The analytical expression follows a piecewise ap-
proach. The analytical expression for the metastable cosine-shaped bump is

$$V(x) = \begin{cases} 0, & -b \le x \le -(a+w) \\ \frac{h}{2}\left(\cos\frac{2\pi(x+2a)}{w} + 1\right), & -(2a+w) \le x \le -2a \\ h + \frac{1}{2}V_0\left(\cos\frac{\pi x}{a} - 1\right), & -2a < x \le -a \\ h + (V_1 - V_0) + \frac{1}{2}V_1\left(\cos\frac{\pi x}{a} - 1\right), & -a < x \le 0 \\ h + (V_1 - V_0) + \frac{1}{2}V_2\left(\cos\frac{\pi x}{a} - 1\right), & 0 < x \le a \\ h + \frac{1}{2}(V_2 - (V_1 - V_0))\left(\cos\frac{\pi x}{a} - 1\right), & a < x \le 2a \\ dh(1-d)\frac{h}{2}\left(\cos\frac{2\pi(x-2a)}{w} + 1\right), & 2a < x \le 2a + w \\ d, & a + w < x \le b \\ \frac{1}{2}k_{\text{harm}}(|x| - b)^2. & |x| > b \end{cases} , \quad (59)$$

where $h$ and $w$ refer to the main barrier height and width, respectively, $d$ is the elevation of
the potential energy at state $B$, and $V_0$, $V_1$ and $V_2$ are the amplitudes of the 3 asymmetric
cosine bumps present. Other parameters carry similar meaning as with the simple cosine
bumps. The the force is again continuous everywhere. For this application we chose $h = 0.8$,
$w = 0.1$, $d = 0.2$, $V_0 = 0.5$, $V_1 = 0.8$, $V_2 = 0.75$ and $b = a + w + 0.1$.

To show robustness of the MSM analysis framework, the interface placement and den-
sity are varied for the metastable bump potential, using diffusive dynamics. A first vari-
ation consists of adding two more interfaces between each existing pair of interfaces to

construct a finer grid (Figure 1, right). A second variation maintained the original number of interfaces, but slightly shifted the inner interfaces such that they no longer coincide with the extrema of the potential. In practice this meant shifting the (innermost) interfaces from 7 equally-spaced interfaces between $\lambda = -0.3$ and $0.3$ to the following positions: $[-0.3, -0.23, -0.14, -0.01, 0.06, 0.24, 0.3]$

The second potential entails a modulated cosine-shaped well centered around $\lambda = 0$, with well-defined stable states $A$ and $B$ on either side, as depicted in Figure 1. Seven equally-spaced interfaces were positioned, and diffusive dynamics were used. This potential is especially interesting as REPPTIS converges significantly faster than RETIS such rugged potentials, where usage of the MSM framework allows accurate estimation of the flux and rate constants.



FIG. 1. The more complex 1D potentials used: the (asymmetric) metastable cosine-shaped bump, with smaller bumps on top (left) and the rugged dip containing many metastable states. The left figure also includes the finer interface grid (indicated in grey) that was used.

[1] A. Ghysels, S. Roet, S. Davoudi, and T. S. van Erp, Physical Review Research **3**, 033068 (2021).

[2] W. Vervust, D. T. Zhang, A. Ghysels, S. Roet, T. S. van Erp, and E. Riccardi, Journal of Computational Chemistry (2024).

[3] L. Verlet, Phys. Rev. **159**, 98 (1967).

12

# 12

# Paper IV (published): Oxygen Storage in Stacked Phospholipid Membranes Under an Oxygen Gradient as a Model for Myelin Sheaths

W. Vervust performed the simulations, and contributed to data analysis and writing of the manuscript.

Check for updates

# Oxygen Storage in Stacked Phospholipid Membranes Under an Oxygen Gradient as a Model for Myelin Sheaths

Wouter Vervust and An Ghysels

**Abstract**

Axons in the brain and peripheral nervous system are enveloped by myelin sheaths, which are composed of stacked membrane bilayers containing large fractions of cholesterol, phospholipids, and glycolipids. The oxygen availability to the nearby oxygen consuming cytochrome $c$ oxidase in the mitochondria is essential for the well-functioning of a cell. By constructing a rate network model based on molecular dynamics simulations, and solving it for steady-state conditions, this work calculates the oxygen storage in stacked membranes under an oxygen gradient. It is found that stacking membranes increases the oxygen storage capacity, indicating that myelin can function as an oxygen reservoir. However, it is found that the storage enhancement levels out for stacks with a large number of bilayers, suggesting why myelin sheaths consist of only 10–300 membranes rather than thousands. The presence of additional water between the stacked bilayers, as seen in cancer cells, is shown to diminish myelin oxygen storage enhancement.

## 1 Introduction

The transport of molecular oxygen throughout the body is primarily realised by blood cells in the cardiovascular system [1]. The final destination of oxygen is the cytochrome $c$ oxidase (COX), a protein embedded in the inner mitochondrial membrane, where oxygen is consumed in the energy conversion cycle. The binding site for oxygen is located at the center of the inner mitochondrial membrane [2]. Passive diffusive transport through several membranes, including the plasma membrane, is responsible for the last steps of oxygen from blood cells to the COX binding site. Cancer cells are known to show hypoxia which can fluctuate across the tumour cell, and tumour oxygenation can enhance radiation therapy of cancer cells [3]. Interestingly, nerve cell axons are surrounded by myelin, which is composed of stacked membrane bilayers containing large fractions of cholesterol, phospholipids, and glycolipids [4, 5]. Understanding the mechanisms of oxygen pathways from capillaries to tissue mitochondria is thus important for completely unraveling the physiological role of oxygen [6].

W. Vervust · A. Ghysels (✉)
IBiTech – Biommeda research group, Ghent University, Ghent, Belgium
e-mail: an.ghysels@ugent.be

301

In this paper, a bottom-up approach will be used to study the role of membranes starting from the molecular scale. Phospholipid membranes form a barrier for oxygen, but the oxygen solubility in the phospholipid tail region is higher than in water, such that a high concentration of oxygen is found in the hydrophobic matrix [7, 8]. This high partitioning toward the membrane suggests that the omnipresence of phospholipid membranes provides the cell with an oxygen reservoir. Therefore, this paper wants to establish the view that a phospholipid membrane has a dual function: slowing oxygen transport, while also storing oxygen. Moreover, the inhomogeneity and anisotropy of the membrane facilitate oxygen transport in the interleaflet region of a lipid bilayer, compared to the transport orthogonal to the membrane surface [9]. The interleaflet space of bilayers can thus form a connected network of efficient diffusive pathways throughout the cell. The timely delivery of oxygen is especially important in the brain and the peripheral nervous system. This paper will in addition investigate how stacked membranes can serve as an oxygen reservoir in the nervous system, given a certain partial oxygen pressure.

In the methods section, kinetic modelling based on the Smoluchowski equation is introduced. It is explained how kinetic parameters were extracted from previous molecular dynamics (MD) simulations. The kinetic parameters will be used to predict the steady-state oxygen distribution given an oxygen gradient. In the results and discussion section, the enhancement in oxygen storage due to the stacking will be computed, where the number of stacked bilayers is varied. The conclusion section summarises our findings.

## 2    Methodology

**Smoluchowski Equation for Transport**  Due to the membrane's inhomogeneity along the membrane's normal ($z$-coordinate, Fig. 1), the oxygen concentration $c$ is $z$-dependent. It relates to the free energy profile $F(z) = -k_B T \ln c(z)$, with $k_B$ the Boltzmann constant and $T$ the temperature.

Diffusive transport through an inhomogeneous medium is governed by the Smoluchowski equation, which is determined by the free energy $F(z)$ and diffusion $D(z)$ profiles [10]. Discretisation in $N$ bins along $z$ transforms the Smoluchowski equation into a set of linearly coupled equations, which can be summarised by a rate model with rate matrix $R$

$$\frac{\partial c_i(t)}{\partial t} = \sum_{N}^{j=1} R_{ij} c_j(t), \tag{1}$$

where $c_i(t)$ is the concentration in bin $i$ at time $t$, and where the kinetic parameters in the rate matrix $R$ are uniquely determined by the $F(z)$ and $D(z)$ profiles [10]. In this work, we will solve the rate equations with different profiles, representing different stacked bilayers.

**Construction of the Rate Model from Molecular Simulations**  Bayesian analysis (BA) of MD trajectories at the atomic scale can be used to obtain the $F$ and $D$ profiles – and thus the rate matrix $R$ – of a membrane [9–11]. In practice, the transitions observed in conventional MD trajectories are compared to the transitions expected by the rate model with a proposed rate matrix $R$. A Monte Carlo procedure is then used to extract the rate matrix with the highest likelihood of producing the observed MD trajectories.

Previously published MD simulations and their $F$ and $D$ profiles are used in this work to investigate the storage capacity of stacked bilayers. The considered membranes are a homogeneous 1-palmitoyl-2-oleoyl-sn-glycero-3-phosphocholine (POPC) bilayer (see Fig. 1) and a model system for inner mitochondrial membrane (labeled MITO) from Ref. [9]. MD simulations were run with the CHARMM36 force field at 310 K with a simulation box containing 72 lipid molecules, as depicted in Fig. 1a [12, 13]. The $z$-axis was discretised into 100 bins ( $z \oplus 0.68^-$ ), and BA was performed on the oxygen trajectories, resulting in the $F$ and $D$ profiles shown in Fig. 1b, c.

In a next step, the single-membrane profiles are extended to artificially construct the $F$ and $D$ profiles of $n$ stacked membranes by concatenat-
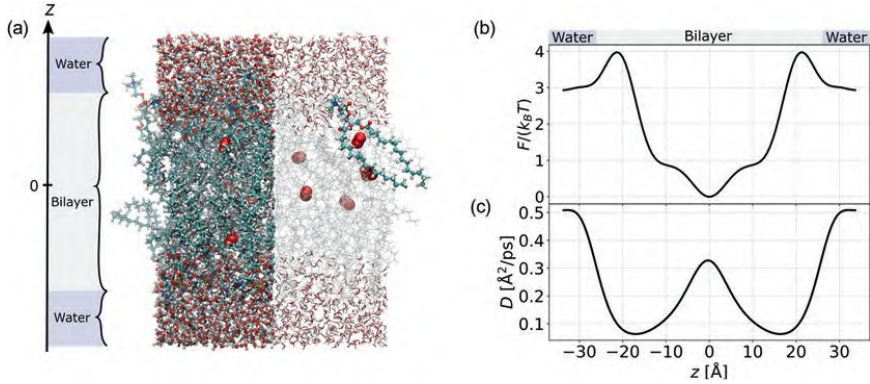
**Fig. 1** MD simulations of the (POPC) membrane. (**a**) Simulation box with indication of the $z$-axis. For clarity, water and phospholipid molecules are shaded in the right part of the box with one high-lighted POPC molecule (ball-and-stick) and red oxygen molecules (van der Waals spheres). (**b**) Free energy $F(z)$ in units $k_B T$. (**c**) Diffusivity $D(z)$. $F(z)$ and $D(z)$ were obtained from Bayesian analysis of the MD trajectories

ing the membrane part of the profiles $n$ times, to imitate myelin sheaths (Fig. 2a). The following three configurations are compared (Fig. 2b):

1. A water layer without any membrane (Fig. 2b top), representing oxygen diffusion to the mitochondria without shielding of membranes.
2. A single bilayer and water (Fig. 2b centre).
3. A series of $n$ stacked bilayers (Fig. 2b bottom), representing myelin sheaths.

The oxygen storage $S$ is measured under an oxygen gradient (see below), which gives $S(0)$, $S(1)$ and $S(n)$ for the setups with zero, one, and $n$ bilayers, respectively. The amount of water is varied in setup 1 and setup 2 to match the size of setup 3, as is depicted in Figs. 2b. This allows for comparison between equally sized configurations in the results and discussion section. The stacking periodicity $d$ can be increased to simulate a more disorganised structure with a larger aqueous phase between the membrane layers.

**Steady-State Solution Under Oxygen Gradient** The oxygen storage $S$ in the membranes is measured under an oxygen gradient, where one side (to the left in Fig. 2) has a partial oxygen pressure and the other side (to the right of Fig. 2) has zero partial oxygen pressure, referring to full depletion of oxygen by COX. The concentration profile is assumed to have reached the steady-state regime, where the concentration no longer changes over time, but a net oxygen flux might still occur.

We solve the rate model (Eq. 1) for the steady-state solution, i.e., the right hand side of Eq. 1 vanishes, with bin 1 (to the left) at fixed concentration $c_1 = c^*$ and bin $N$ (to the right) at fixed concentration $c_N = 0$, for all times $t$. This is done by constructing the vector $v$ of length $N - 2$, which has all elements equal to zero except for its first element $v_1 = -R_{2,1}$. The rate matrix is stripped of its first and $N^{th}$ rows and columns, to form the cropped rate matrix $R'$ [14]. The steady-state concentration $c_i^{ss}$ in each bin $i$ is then obtained by matrix multiplication of $R'$ and the vector $v$, and scaling with the concentration $c^*$,

$$c_i^{ss} = \begin{cases} c^*, & i = 1 \\ c^* \sum_{j=2}^{N-1} \left[ R'^{-1} \right]_{i-1, j-1} v_j, & i = 2, \ldots, N-1 \\ 0, & i = N \end{cases} \tag{2}$$
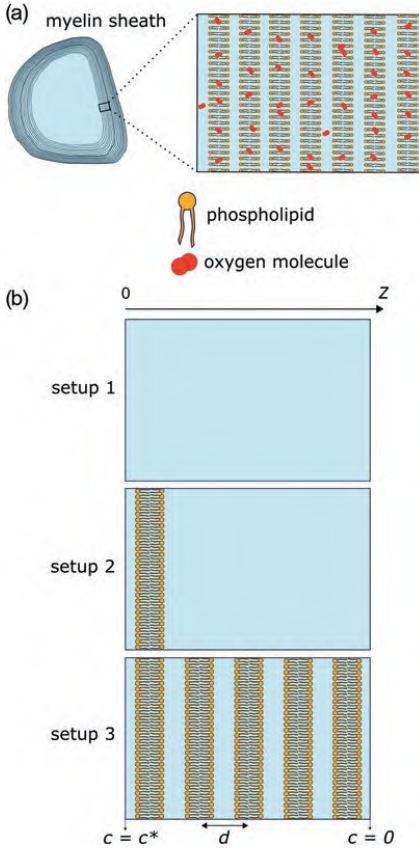
**Fig. 2** (**a**) Myelin sheaths wrapped around a neural cell consist of stacked membranes [5]. (**b**) Considered configurations with steady-state boundary conditions: $c = c^*$ to the left and $c = 0$ to the right. Setup 1 is pure water, setup 2 has one bilayer, setup 3 has $n$ bilayers (depicted is $n = 5$). Stacking periodicity $d$ is indicated with arrow

where $R'^{-1}$ is the matrix inverse of $R'$. Finally, the storage capacity $S$ corresponding to a concentration profile $c(z)$ is computed by integrating over the concentration profile and multiplying with the membrane surface area $A$,

$$S = A \int c(z) dz \approx A \sum_i c_i \ z \tag{3}$$

where in the second equation, the integral was approximated by the discretisation with bin width $\Delta z$.

## 3 Results and Discussion

First, we calculated the concentration profile under an oxygen gradient (Eq. 2) in the setups without bilayer, with one bilayer, and with $n$ stacked bilayers. The steady-state profiles for the POPC membrane are shown for $n = 5$ in Fig. 3a. The setup of pure water gives a simple linear slope (black line), while the presence of a membrane clearly introduces a high peak in oxygen concentration (blue dashed line). Increasing the number of stacked membranes creates additional peaks (green line), whose height drops linearly with each additional peak.

The stacking enhancement of adding bilayers is quantified by $E(n) = S(n)/S(1)$, which is the ratio between storage in $n$ stacked bilayers versus storage in 1 bilayer (setup 2 and setup 3 in Fig. 2). The membrane area $A$ (Eq. 3) drops out of this ratio, and $E(n)$ has no unit. Fig. 3b depicts the stacking enhancement $E(n)$ for the POPC and MITO membranes, which shows that the stacking enhancement increases with the number $n$ of stacked membranes. This indeed suggests that myelin sheaths may function as a reservoir for oxygen storage in the nervous system, making oxygen readily available to COX in the nearby mitochondria.

Interestingly, stacked bilayers not only have more peaks in Fig. 3a, but the first peak of the stacked bilayers (green) is also higher compared to the case of a single bilayer (blue dashed). This again favours the enhancement in stacked membranes compared to a single bilayer. The higher peak is explained by the fact that the nearby second peak also provides a source of oxygen.

Moreover, the slope of the curve in Fig. 3b becomes less steep with increasing $n$, meaning that stacking an additional membrane onto a thick layer of stacked membranes has less impact. In other words, stacking a 51st membrane on top of 50 stacked membranes has less impact on the stacking enhancement than stacking an 11th membrane on top of 10 stacked membranes. Therefore, it seems plausible that, for a certain stacking number, a trade-off may be reached between stacking enhancement and the extent of the physical space occupied by the membranes. Indeed, in electron microscopy images, myelin

12.   Paper IV (published): Oxygen Storage in Stacked
Phospholipid Membranes Under an Oxygen Gradient as a
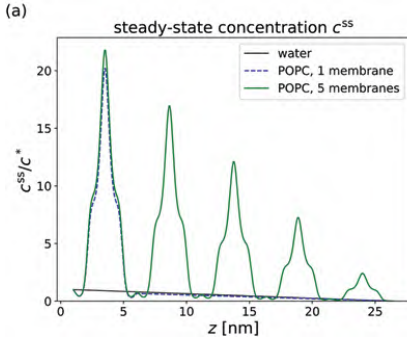Model for Myelin Sheaths

**Fig. 3** (**a**) Steady-state concentration profile under oxygen gradient (in units $c^*$) with $c = c^*$ to the left and $c = 0$ to the right, for $n = 5$. Setup 1 with only water (black), setup 2 with 1 POPC membrane (blue dashed) and setup 3 with $n = 5$ stacked POPC membranes (green), as shown in
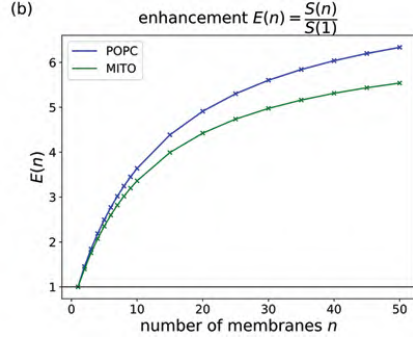
Fig. 2b. (**b**) Stacking enhancement $E(n) = S(n)/S(1)$ is the ratio of storage capacity between n stacked bilayers and one bilayer, under a steady-state oxygen gradient over the membrane(s). $S(1)$ and $S(n)$ correspond to setup 2 and setup 3 in Fig. 2, respectively

sheaths often consist of about 10–300 bilayers but not thousands [5].

The storage efficiency also depends on the intrinsic storage capacity of a single bilayer. The ratio $S(1)/S(0)$ is a measure for the content in a single bilayer compared to a slab of pure water with the same dimension. For POPC, this ratio amounts to $S(1)/S(0) \approx 8$ in the steady-state regime using a membrane thickness of approximately 5.2 nm. Because of the linearity of the rate model and the symmetry of the considered setups, we can conclude that the latter ratio is equivalent to the ratio of storage in a membrane versus storage in the water phase in equilibrium. This is also known as the membrane partitioning constant $K_m$, of which the study was initiated over a century ago with Overton's rule. In short, $S(1)/S(0) = K_m = c_{mem}/c_{wat}$, with $c_{mem}$ and $c_{wat}$ the equilibrium oxygen concentration in the membrane and water phase, respectively.

Concretely, the POPC membrane stores about 8 times more oxygen molecules than can be stored in pure water. For example, let us assume a partial oxygen pressure of ~10 mmHg, which lies in the reported 10–34 mmHg range for brain tissue [15]. At atmospheric pressure and body temperature, the partial oxygen pressure is 159 mmHg, and the solubility of oxygen in water is approximately 7.187 mg/L [16]. Therefore, the

10 mmHg partial oxygen pressure corresponds to an oxygen concentration of $(10/159) \cdot 7.187$ mg/L = 0.45 mg/L. Next, the steady-state profile over 10 stacked POPC membranes under this oxygen gradient ($c^* = 0.45$ mg/L to the left of the stack, while $c = 0$ to the right of the stack) is computed with Eq. 2, and similarly for pure water. Finally, the storage per area ($S/A$) follows from Eq. 3, resulting in about $1791/\mu m^2$ oxygen molecules available in the membrane stack, compared to $226/\mu m^2$ oxygen molecules in a similarly sized slab of water.

Next, it is tested how the bilayer spacing influences the enhancement provided by the stacked membranes. The pathophysiology of cancer cells may show abnormal microvasculature and defective microcirculatory function, such as an increase in interstitial fluid [17, 18]. Moreover, abnormal myelination can also occur following trauma to a central nervous system tract, where myelin sheath swelling and decompaction have been observed, even in myelin that is spatially separated from the primary injury [19, 20]. Here, structural loss in the membrane stacking is modeled by increasing the thickness of the aqueous phase between the stacked bilayers. These additional water layers will decrease the total permeability [21], but they might also affect the storage capacity. Fig. 4 shows the

**Fig. 4** Enhancement $S(n)/S(0)$ in stacked membranes for $n = 10$, as a function of stacking periodicity $d$



enhancement $S(n)/S(0)$ for $n = 10$ stacked bilayers, comparing the storage in 10 bilayers with the storage in pure water, while the stacking periodicity $d$ (indicated in Fig. 2b) is gradually increased. At the left of the figure, the periodicity $d = 5.2$ nm corresponds to 10 stacked POPC membranes without additional water, as was simulated in Fig. 3b. Figure 4 clearly shows that an increase in stacking periodicity $d$ can cause a strong decrease in the storage enhancement by the stacked bilayers, indicating that additional water between the myelin layers affects the storage efficiency in a negative way. This can be understood by comparing the contributions to the total storage by the membrane and by the interstitial water. The oxygen storage in water becomes larger for increasing $d$, simply because there is more interstitial water. In contrast, the oxygen storage in the membranes is fairly independent of $d$. Thus, for increasing periodicity $d$, the relative contribution of the bilayers to the total oxygen storage becomes smaller, and consequently a highly distorted membrane stack will gradually lose its function as a compact oxygen reservoir.

## 4    Conclusion

A phospholipid membrane has a higher solubility for molecular oxygen compared to pure water. By calculating the steady-state concentration over a membrane under an oxygen gradient, we have shown that the oxygen storage capacity $S$ can be even more enhanced by stacking multiple bilayers, as is the case in myelin sheaths of axons in the nervous system. The stacking enhancement levels out when many membranes are stacked, meaning that beyond a certain number of stacked membranes, stacking even more would have limited efficiency. This could explain why myelin sheaths appear to consist of 10–300 stacked bilayers but not thousands. Finally, the enhancement is negatively affected by an increase of the aqueous phase between the bilayers, which can for instance occur in cancer cells with more disordered stacking geometry.

## References

1. Popel AS (1989) Theory of oxygen transport to tissue. Crit Rev Biomed Eng 17:257–321
2. Tsukihara T, Aoyama H, Yamashita E et al (1995) Structures of metal sites of oxidized bovine heart cytochrome c oxidase at 2.8 A. Science 269:1069–1074
3. Colliez F, Gallez B, Jordan BF (2017) Assessing tumor oxygenation for predicting outcome in radiation oncology: a review of studies correlating tumor hypoxic status and outcome in the preclinical and clinical settings. Front Oncol 7:10

12.   Paper IV (published): Oxygen Storage in Stacked
Phospholipid Membranes Under an Oxygen Gradient as a
Model for Myelin Sheaths

4. Morell P, Quarles RH (1999) Characteristic composition of myelin. Basic Neurochem Mol Cell Med Asp 6:69–94

5. Arroyo EJ, Scherer SS (2000) On the molecular architecture of myelinated fibers. Histochem Cell Biol 113:1–18

6. Pias SC (2021) How does oxygen diffuse from capillaries to tissue mitochondria? Barriers and pathways. J Physiol 599:1769–1782

7. Dotson RJ, Smith CR, Bueche K et al (2017) Influence of cholesterol on the oxygen permeability of membranes: insight from atomistic simulations. Biophys J 112:2336–2347

8. Ghysels A, Krämer A, Venable RM et al (2019) Permeability of membranes in the liquid ordered and liquid disordered phases. Nat Commun 10:1–12

9. Ghysels A, Venable RM, Pastor RW, Hummer G (2017) Position-dependent diffusion tensors in anisotropic media from simulation: oxygen transport in and through membranes. J Chem Theory Comput 13:2962–2976

10. Hummer G (2005) Position-dependent diffusion coefficients and free energies from Bayesian analysis of equilibrium and replica molecular dynamics simulations. New J Phys 7

11. Krämer A, Ghysels A, Wang E et al (2020) Membrane permeability of small molecules from unbiased molecular dynamics simulations. J Chem Phys 153:124107

12. Brooks BR, Brooks CL III, Mackerell AD Jr et al (2009) CHARMM: the biomolecular simulation program. J Comput Chem 30:1545–1614

13. Klauda JB, Venable RM, Freites JA et al (2010) Update of the CHARMM all-atom additive force field for lipids: validation on six lipid types. J Phys Chem B 114:7830–7843

14. Vervust W, Ghysels A (2022) An electric RC circuit model to describe oxygen storage and transport in myelin sheaths. To be submitted

15. Mulkey DK, Henderson RA III, Olson JE et al (2001) Oxygen measurements in brain stem slices exposed to normobaric hyperoxia and hyperbaric oxygen. J Appl Physiol 90:1887–1899

16. Clever LH, Battino R, Miyamoto H et al (2014) IUPAC-NIST solubility data series. 103. Oxygen and ozone in water, aqueous solutions, and organic liquids (supplement to solubility data series volume 7). J Phys Chem Ref Data Monogr 43:33102

17. Molls M, Anscher MS, Nieder C, Vaupel P (2009) The impact of tumor biology on cancer treatment and multidisciplinary strategies. Springer

18. Vaupel P (2006) Abnormal microvasculature and defective microcirculatory function in solid tumors. In: Vascular targeted therapies in oncology. Wiley, pp 9–29

19. Payne SC, Bartlett CA, Harvey AR et al (2011) Chronic swelling and abnormal myelination during secondary degeneration after partial injury to a central nervous system tract. J Neurotrauma 28:1077–1088

20. Payne SC, Bartlett CA, Harvey AR et al (2012) Myelin sheath decompaction, axon swelling, and functional loss during chronic secondary degeneration in rat optic nerve. Invest Ophthalmol Vis Sci 53:6093–6101

21. Davoudi S, Ghysels A (2021) Sampling efficiency of the counting method for permeability calculations estimated with the inhomogeneous solubility–diffusion model. J Chem Phys 154:54106

# 13

# PAPER V (SUBMITTED): MYELIN SHEATHS CAN ACT AS COMPACT TEMPORARY OXYGEN STORAGE UNITS AS MODELED BY AN ELECTRICAL RC CIRCUIT MODEL

# Myelin sheaths can act as compact temporary oxygen storage units as modeled by an electrical RC circuit model

**Wouter Vervust**[1], **Katja Witschas**[2], **Luc Leybaert**[2], and **An Ghysels**[1,*]

[1]IBiTech - BioMMedA group, Ghent University, Belgium
[2]Physiology Group, Department of Basic and Applied Medical Sciences, Ghent University, Belgium
[*]an.ghysels@ugent.be

## ABSTRACT

Oxygen is a crucial component in cellular energy metabolism, particularly in the brain where neurons consume the majority of energy produced via mitochondrial oxidative phosphorylation. Phospholipid membranes have been found to store and transport oxygen efficiently in their hydrophobic core. This work investigates the kinetics of this storage, not only in a single membrane but also in membrane stacks as those found in myelinated axons in the nervous system. Using a diffusive model derived from molecular dynamics simulations, it is first demonstrate that oxygen storage within a phospholipid bilayer follows first-order kinetics. In consequence, we show how oxygen loading and unloading in a membrane is effectively modeled by an intuitive RC (resistor-capacitor) circuit analogy with a characteristic RC time constant. Next, oxygen transport through myelin, comprising multiple bilayers, could be investigated by building a ladder network of RC circuits. Both the resistance (to oxygen transport) and the capacitance (for oxygen storage) scale linearly with the number of bilayers. Moreover, the characteristic time constant for oxygen storage scales quadratically with the myelin thickness, for instance enhancing the characteristic time constant from 30 ns for one 1-palmitoyl-2-oleoyl-sn-glycero-3-phosphocholine (POPC) bilayer to 506 $\mu$s for 200 POPC bilayers. This enhancement gives myelin a buffering role: myelin sheaths act as spatially compact oxygen containers whose slower kinetics will dampen sudden sudden oxygen changes. Finally, oxygen transport from capillary to axonal mitochondria is modeled during increased oxygen consumption rates (OCRs) associated with neuronal activity. The model predicts that increased myelination results in longer sustainment of increased oxygen demand, supporting the idea that myelin sheaths may act as a buffer to oxygen fluctuations. The inability of considered configurations to sustain the increased OCRs for periods longer than a few hundred milliseconds hints the functional aspect of the dominant vascular response in restoring oxygen homeostasis following neuronal activation.

## 1 Introduction

The brain covers over 20 % of total oxygen metabolism, and it is estimated that neurons consume 75 % to 80 % of energy produced in the brain[1,2]. Central to this energy production is the mitochondrial aerobic synthesis of adenosine triphosphate (ATP) via oxidative phosphorylation (OxPhos). Oxygen is the terminal electron acceptor in the electron transport chain (ETC), and therefore essential in providing the proton gradient necessary for ATP synthesis[3].

The mechanism of oxidative metabolism is not yet fully understood, and the classical picture of mitochondrial monopoly over CNS energy production has been challenged. Indications for extra-mitochondrial aerobic ATP synthesis via OxPhos within the myelin sheath have been presented[4–9]. Oligodendrocytes and Schwann cells have also been reported to aid axon metabolism, with abnormal oligodendrocyte and Schwann cell metabolism leading to axon degeneration[10–13].

Nevertheless, cytochrome $c$ oxidase (COX) remains the most crucial enzyme for molecular oxygen in the brain. COX's high affinity for oxygen ensures undiminished activity of OxPhos even at very low $P_{O_2}$[14–16]. However, a significant drop in oxygen supply can halt OxPhos, drastically reducing ATP production and disrupting ion gradients.

As oxygen consumption is central to energy production, it can be used to probe neuronal activity. The blood oxygen level dependent (BOLD) signal detected in functional magnetic resonance imaging (fMRI) often displays an 'initial dip' in gray matter, which has been linked to a quick increase in the cerebral metabolic rate of oxygen (CMRO$_2$) prior to the cerebral blood flow (CBF) increase[17,18]. The initial dip in white matter is likely prolonged compared to gray matter[19,20]. The vascular response (CBF increase) occurs $\sim$500 ms after the neural stimulus, and quickly thereafter dominates the increase of CMRO$_2$[21,22]. Directly probing $P_{O_2}$ in capillaries using two-photon lifetime microscopy (2PLM) also displays the initial dip in tissue $P_{O_2}$ following neuronal activation[23,24]. This dip is observed within $\sim$100 ms following synaptic activation, and preceding functional hyperemia by $\sim$1 s[23]. The shape and duration of the dip are, however, unpredictable, as they depend highly on local parameters[23,25]. While other probing techniques exist, there remains a fundamental knowledge gap at subcellular spatial scales <1 $\mu$m. The pathways of oxygen delivery towards COX are complex, consisting of many energetic and diffusive barriers, where the impact of subcellular structures on oxygen transport phenomena remains largely unknown. This work investigates the impact of axon myelination on oxygen kinetics, contributing to closing this subcellular knowledge gap.

1

Myelin sheaths, consisting of up to 100 tightly packed phospholipid bilayers, insulate axons and provide the necessary speed-up for electric signals via saltatory conduction of action potentials[26]. Along the axon, myelinated regions (internodes) alternate unmyelinated regions (nodes of Ranvier), distributed quasi equidistantly along the axon[27,28]. Total oxygen solubility is 3 to 5 times larger in phospholipid bilayers than in water, reaching upwards of 10 times in the hydrophobic mid-plane of the bilayer[29–33], suggesting that myelin may serve as an oxygen reservoir at close proximity to the axonal mitochondria[34]. Even though the tightly packed polar phospholipid headgroups pose both an energetic as well as a diffusive barrier, membrane oxygen permeability is still relatively high as predicted by Overton's rule (high overall oxygen solubility) and as evidenced by molecular dynamics studies[35–38]. Both oxygen solubility and diffusivity in water decrease significantly with increasing protein and solute content, decreasing the diffusion barrier depth, and ultimately enforcing a dense capillary distribution[39]. Sufficient oxygenation of brain tissue during gamma oscillations was estimated to be possible only within a radius of 30 to 40 μm from the capillary, which lies closely to capillary distances of 40 to 70 μm in humans and rodents[40].

In this work, oxygen transport through myelin is investigated. A kinetic model is built from MD simulation data, and it is shown that oxygen depletion and filling of phospholipid bilayers is governed by first-order kinetics. Next, a simple RC circuit model is constructed that intuitively describes oxygen transport using 2 membrane parameters: the membrane resistance (to oxygen permeability) and the membrane capacitance (to store oxygen). Viewing a myelin sheath as a stack of phospholipid bilayers, oxygen transport through myelin is then modeled as a ladder circuit of bilayer elements. It is shown that the total resistance and capacitance of myelin increase linearly with the number of bilayers, whereas the characteristic time of oxygen storage increases quadratically. Oxygen transport from capillaries to axonal mitochondria is then modeled by extending the myelin model with (extracellular and cytosolic) solvent phases with variable oxygen consumption boundary conditions. It is shown that myelination increases the duration for which an internodal axon segment can support increased oxygen demands following neuronal activity. Furthermore, our model predicts that oxygen demand cannot be met for durations larger than 100 ms, supporting the functional aspect of the hyperemic CBF response that quickly follows thereafter.

## 2 Results

### Oxygen permeation through a single bilayer is a first-order kinetics process

Assuming diffusive dynamics, the Smoluchowski model describes the oxygen distribution in a lipid bilayer of thickness $h$, governed by the position-dependent free energy $F$ and diffusivity $D$ profiles. Both profiles were extracted for molecular oxygen in a 1-palmitoyl-2-oleoyl-sn-glycero-

3-phosphocholine (POPC) bilayer using Bayesian analysis on equilibrium molecular dynamics (MD) simulations (see Fig. 1A for the simulation box). To examine the time-dependent behavior of oxygen with arbitrary initial concentration in the membrane and various boundary conditions (BCs) in the solvent region, the Smoluchowski model is approximated by a set of rate equations via spatial discretization of the Smoluchowski equation into $N$ bins. This discretized Smoluchowski rate model is characterized by the rate matrix $R$ of size $N \times N$. Additional details regarding the MD simulations, the Smoluchowski model, and the discretized Smoluchowski rate model are provided in the Methods section.

The Smoluchowski rate model is now solved (Eq. 5) for the case of an oxygen-deprived POPC bilayer as initial condition ($c(0,z) = 0$, $0 < z < h$ at $t = 0$) and exposed to a source of constant oxygen concentration in the solvent region at the left-hand leaflet as BCs ($c(t,0) = c_L$, $c(t,h) = 0$, $\forall t > 0$). Fig. 1B shows how oxygen propagates into the bilayer and gradually reaches the steady-state distribution $c^{ss}(z)$. The oxygen concentration quickly becomes largest at the bilayer mid-plane which is expected as oxygen prefers the hydrophobic core over the polar lipid head-group regions. The time-dependency of the oxygen storage $S(t)$, defined as the amount of oxygen molecules in the POPC bilayer per area of the membrane cross section (Eq. 3), is shown in Figs. 1C and in a log plot in Fig. 1D.

Interestingly, $S(t)$ closely adheres to first-order kinetics. The storage converges almost exponentially to the steady-state oxygen storage $S^{ss} = S(t \to \infty)$ under the given BCs. The single exponential approximation $S'(t) = S^{ss}(1 - \exp(-t/\tau))$ is added to Figs. 1C-D for comparison. Here, $\tau$ is the largest time constant associated to the largest non-zero eigenvalue of the rate matrix of the Smoluchowski rate model, and thus describes the slowest time scale in the oxygen loading of the membrane. The exponential approximation $S'(t)$ closely matches $S(t)$, except on small time scales where the faster processes dominate. The effect of these faster processes on $S(t)$ is not visible in Figs. 1C-D as their small contribution to $S(t)$ becomes negligible for $t \gtrsim 2$ ns.

### Building a membrane RC circuit with membrane resistance and membrane capacitance

The remarkable single exponential behavior of the POPC membrane filling with oxygen allows us to draw a parallel with the charging of a capacitor in an electric RC (resistor-capacitor) circuit, as both systems exhibit first-order kinetics.

Electric circuit analogies are often performed for transport phenomena, such as the Hodgkin-Huxley model for transport of charged particles and electric currents in neurons[41]. In relation to oxygen transport specifically, such analogies have been done to analyze facilitated transport in solid membranes with fixed site carriers[42], or to investigate the oxygen diffusion impedance in proton exchange membrane fuel cells[43], or even to model phase transitions of red blood cells to sickle cells[44].

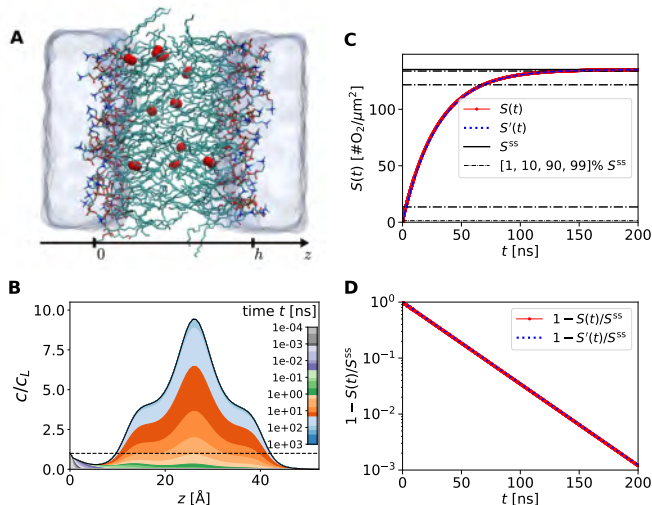Similarly to an RC circuit collecting charge in its capac-

**Figure 1. A**: Snapshot of the MD simulation, containing the POPC lipids (lines), molecular oxygen (red van der Waals spheres), and water molecules (visualized as a transparent surface). Panels **B-D** study a POPC bilayer that is initially oxygen-depleted at $t < 0$. The boundary conditions (BCs) fix the left solvent oxygen concentration to $c_L = 0.45\,\mathrm{mg/L}$ and the right at $0\,\mathrm{mg/L}$ for $t > 0$. **B**: Oxygen concentration profile $c(t,z)/c_L$ in the POPC bilayer evolves towards steady-state concentration $c^{ss}(z)$ (thick black line). Dashed horizontal line denotes $c/c_L = 1$. Panels **C-D**: Oxygen storage $S(t)$ (red) evolves towards the steady-state storage $S^{ss}$ of 135 $O_2$ molecules per 1 μm$^2$ slab of POPC. Black horizontal lines represent fractions of $S^{ss}$. The approximation $S'(t)$ (blue line) is based on the largest time constant $\tau$ of the discretized Smoluchowski model. **D**: The time-dependence of $S(t)$ is well approximated by a single exponential, as is evident from comparing $1 - S(t)/S^{ss}$ with $1 - S'(t)/S^{ss} \equiv \exp(-t/\tau)$ on a logarithmic plot.

itor by conducting current through a resistor, a membrane can store oxygen in its hydrophobic core by letting oxygen permeate through the leaflets. A clear connection between membrane permeation properties and electric current properties will now be established, while implications will be brought up in the Discussion section. Two components are central to this analogy: the membrane resistance ($R_M$) and the membrane capacitance ($C$). These components characterize oxygen behavior in a membrane and are fundamental for constructing an electric circuit that accurately models oxygen kinetics within the POPC bilayer.

First, consider the resistance to oxygen transport through the membrane. The net oxygen flux $J$ through a membrane is proportional to the difference in oxygen concentrations $\Delta c$ in the left ($c_L$) and right ($c_R$) solvent regions. The proportionality constant $R_M$ is the membrane resistance and is the reciprocal of the well-known membrane permeability $P$. This proportionality relation, $\Delta c = J R_M = J/P$, corresponds to Ohm's law for electric resistors, where oxygen concentration $c$ is analogous to voltage, and oxygen flux $J$ is analogous to electric current. To account for differences in lipid content

between bilayer leaflets, the membrane resistance is divided into a series of left ($R_L$) and right ($R_R$) leaflet resistances, such that $R_M = R_L + R_R$.

Second, consider the membrane's capability to store oxygen in its hydrophobic core. At biological oxygen concentrations, there are no oxygen crowding effects, and therefore the oxygen concentration $c(z,t)$ is proportional to both the BCs at $c_L$ and $c_R$ (linearity of Smoluchowski equation). In other words, increasing either the extracellular or intracellular oxygen concentration increases the amount of oxygen in the bilayer. Consequently, at equilibrium, where the net oxygen flux and concentration difference are both zero, the oxygen stored in the membrane $S^{eq}$ (per unit of cross area) is proportional to the solvent oxygen concentration $c_L = c_R = c_{ref}$. We now define this proportionality constant as the membrane capacitance $C$, i.e. $S^{eq} = C c_{ref}$. In this context, the oxygen storage $S^{eq}$ is comparable to the electric charge stored in an electric capacitor. However, unlike an electric capacitor where charge depends solely on the voltage difference, the oxygen membrane capacitance depends on both $c_R$ and $c_L$ rather than just the concentration difference $\Delta c$. Rather than defining two

| Electric | | Membrane | |
|---|---|---|---|
| Property | Unit | Property | Unit |
| voltage $V$ | Volt | concentration $c$ | mol/m$^3$ |
| current $I$ | Ampere | flux $J$ | mol/(m$^2$s) |
| charge $Q$ | Coulomb | storage $S$ | mol/m$^2$ |
| capacitance $C$ | Farad | capacitance $C$ | m |
| resistance $R$ | Ohm | resistance $R$ | s/m |
| Capacitive laws | | | |
| $dQ(t)/dt = I(t)$ | | $dS(t)/dt = J_{\text{eff}}(t)$ | |
| $dV(t)/dt = I(t)/C$ | | $dc_{\text{eff}}/dt = J_{\text{eff}}(t)/C$ | |
| $Q(t) = CV(t)$ | | $S(t) = Cc_{\text{eff}}(t)$ | |
| Resistive laws | | | |
| | | $c_L(t) - c_{\text{eff}}(t) = R_L J_L(t)$ | |
| $V(t) = RI(t)$ | | $c_R(t) - c_{\text{eff}}(t) = R_R J_R(t)$ | |

**Table 1.** Equivalence between electric circuit properties and membrane permeation properties. Their respective units, capacitive laws, and resistive laws are given. The circuit is depicted in Fig. 2.



**Figure 2.** Membrane RC circuit modeling time-dependent oxygen kinetics in a POPC bilayer. The solvent regions (blue frames) are represented by the voltage sources $c_L(t)$ and $c_R(t)$. The membrane (red frame) is given by the junction of the two leaflet resistances $R_L$ and $R_R$, and the membrane capacitance $C$. Also indicated are the currents $J_L(t)$, $J_R(t)$, $J_{\text{eff}}(t)$, and the voltage $c_{\text{eff}}(t)$. After solving the circuit (Eq. 9), the net flux $J_{\text{eff}}(t)$ through the membrane is obtained as $J_{\text{eff}}(t) = J_R(t) - J_L(t)$, and the oxygen storage is obtained as $S(t) = Cc_{\text{eff}}(t)$.

proportionality constants, we opt to define $C$ with the equilibrium storage where $c_L$ and $c_R$ are equal. This means that an increase in the oxygen concentration in either the extracellular or intracellular solvent regions will result in an increase in oxygen storage in the bilayer, even if $\Delta c$ remains constant. The analogy between electric and membrane properties is summarized in Table 1 with their respective units.

The circuit in Fig. 2 models the time-dependent oxygen kinetics in a POPC membrane submerged in solvent. The membrane consists of the membrane capacitance $C$ and the two leaflet resistances $R_L$ and $R_R$, while the solvent regions are given by the time-dependent voltage sources $c_L(t)$ and $c_R(t)$. The currents $J_L(t)$ and $J_R(t)$ represent the oxygen influx at the left leaflet-solvent interface and the oxygen efflux at the right leaflet-solvent interface, respectively. The difference of these fluxes gives the total (or effective) oxygen flux from the membrane to the solvent $J_{\text{eff}}(t) = J_R(t) - J_L(t)$. The oxygen concentrations are treated as voltage potentials in the circuit, and the oxygen storage $S(t)$ in the membrane is derived from the 'voltage' over the capacitor $C$ as $S(t) = Cc_{\text{eff}}(t)$. Thus, the concentration $c_{\text{eff}}(t)$ is not a physical concentration in the membrane, but rather an 'effective concentration' for which the capacitance formula can be used.

This circuit correctly models oxygen behavior at (1) equilibrium, (2) steady-state, and (3) while (dis)charging the membrane. Firstly, at equilibrium, the oxygen concentrations at both sides of the bilayer are equal ($c_{\text{eff}} = c_L = c_R = c_{\text{ref}}$) such that the equilibrium storage $S^{\text{Eq}} = c_{\text{ref}}C$ is recovered. Secondly, at steady-state ($c_L \neq c_R$), the capacitor becomes an open circuit such that the impedance between the left and right voltage sources is equal to the sum of the leaflet resistors $R_L + R_R = R_M$. As such, the permeability is recovered via
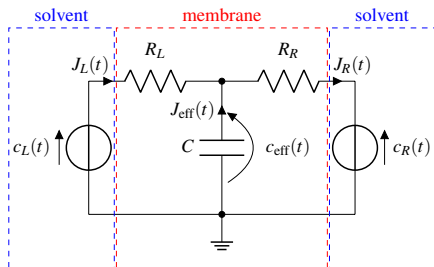
Ohm's law $J_R = J_L = \Delta c/R_M = \Delta c P$. Thirdly, the time constant of the circuit is given by $\tau_{\text{RC}} = C \frac{R_L R_R}{R_L + R_R}$, which is a very good approximation of the characteristic time scale $\tau$ of membrane loading. For POPC, which has equal leaflet resistances ($R_L = R_R = R_M/2$), the RC circuit time constant simplifies to $\tau_{\text{RC}} = R_M C/4 = 31.263$ ns, which is only a 4% overestimation of the characteristic time constant $\tau = 29.952$ ns of the discretized Smoluchowski rate model.

The deviation between $\tau$ and $\tau_{\text{RC}}$ can be understood by looking at their origin. On one hand, $\tau_{\text{RC}}$ follows from an RC circuit with the true membrane resistance $R_M$ (Eq. 2) and true membrane capacitance $C$ (Eq. 6). On the other hand, the characteristic time scale $\tau$ is the largest time scale of the discretized Smoluchowski rate model, associated to the highest non-zero eigenvalue. In an algebraic formulation, one could say that $\tau$ is the time constant associated with the storage part included in the largest eigenvector. This implies that $\tau$ does not cover the full complex dynamics of membrane loading or depletion, as $\tau$ neglects the existence of all other fast time scales. Nevertheless, as discussed above, $\tau$ provides a very fair approximation (Fig. 1C-D), and the numerical comparison between $\tau$ and $\tau_{\text{RC}}$ for POPC indicates that the RC circuit in turn reproduces this time scale well with 4% accuracy. Overall, we conclude that Fig. 2 is the electrical circuit representing oxygen kinetics for storage and loading/depletion for a bilayer.

## Ladder circuit to model a myelin sheath

A myelin sheath can be modeled as a stack of $M$ bilayers, where oxygen transport is simulated using two techniques. Firstly, a discretized Smoluchowski rate model for the $M$-

**4**/17

177

stack is constructed by appending the $F$ ($D$) profile of POPC $M$ times to obtain the free energy (diffusivity) profile of the $M$-stack. The rate matrix associated to these profiles is then constructed, the constant oxygen concentration BCs are applied, and the time-dependence of the oxygen concentration in the stack follows from solving the rate model (Eq. 5). This technique serves as the benchmark.

Secondly, the sheath is treated as a series of $M$ single bilayers, where we model each bilayer by an RC circuit (red frame in Fig. 2). The sheath circuit (Fig. 3B) is constructed by connecting $M$ membrane circuits, which is then connected to left and right oxygen sources in the solvent, $c_L(t)$ and $c_R(t)$, respectively. By solving the resulting ladder circuit for its 'voltages' $c_{\text{eff}}^{(i)}(t)$, the oxygen storage and oxygen flux through the stack can be computed for different initial conditions and BCs of the oxygen sources, as shown in the Methods section (Eq. 11).

The evolution of $S(t)$ in a myelin sheath consisting of $M = 5$ POPC bilayers is illustrated in Fig. 3C-D. A charging scenario similar to that of the single POPC bilayer is considered: the membrane stack is initially oxygen-depleted ($S(t) = 0$, $t < 0$), after which $c_L$ is suddenly set to $c_L$ whereas $c_R$ remains zero. The evolution of the oxygen concentration profile $c(z,t)$ in the myelin sheath is shown in Fig. 3C for the rate matrix technique (Eq. 5). The left-most bilayer takes most of the initial oxygen at the shorter time scales, after which the oxygen permeates to the other bilayers at the longer time scales. Consider the oxygen storage $S_i(t)$ in the $i$-th bilayer, counting from left to right, with its steady-state storage $S_i^{\text{ss}} = S_i(t \to \infty)$. The normalized storage $S_i(t)/S_i^{\text{ss}}$ in Fig. 3D shows that the bilayer that is furthest from the oxygen source (i.e. for bilayer 5) depends most on the large-timescale contributions, since it reaches its $S_5^{\text{ss}}$ value the latest. Fig 3D shows the normalized storages calculated with the two techniques, comparing the simplified RC ladder model (colors, Eq. 5) to the discretized Smoluchowski rate model (black, reference, Eq. 11). The difference between the techniques is hardly visible, which indicates that the reference is approximated well by the ladder network. Oxygen transport in a myelin sheath can thus be considered as a ladder of RC circuits where the 'oxygen capacitors' gradually get 'charged' with oxygen molecules, which permeate with a permeability $P$ through each bilayer.

### Dependence on thickness of myelin sheaths

The total resistance (capacitance) of an $M$-stack myelin sheath will be sum of the individual bilayer resistances (capacitances). Both the total amount of oxygen storage and the oxygen response time also depend on the size of the myelin sheath. Fig. 3E shows the total amount of oxygen stored $S^{\text{eq}}$ in myelin sheaths containing $M = 1$ to 200 POPC bilayers in equilibrium with solvent regions with oxygen concentration $c_{\text{ref}}$. The Smoluchowski technique (black) based on the discretized free energy (Eq. 4) is given as the reference and shows that the storage increases linearly with the number of POPC bilayers $M$. The approximate technique of the ladder circuit prediction

(red, Eq. 9) again agrees remarkably well. At equilibrium, all effective potentials $c_{\text{ref},i}$ are equal to $c_{\text{ref}}$, and the total storage is $S^{\text{eq}} = c_{\text{ref}} \sum_{i=1}^{M} C$. In case of a stack of bilayers with identical capacitance $C$, the storage is thus linear in the number $M$ of bilayers, $S^{\text{eq}} = c_{\text{ref}} MC$.

The larger the myelin stack, the slower the total oxygen storage responds to changes in external oxygen concentrations. This can be understood by looking at the largest time constant $\tau$ of the ladder circuit (the dominant non-zero eigenvalue of the associated rate matrix $A$, see Eq. 11). The (almost) quadratic increase of $\tau$ with the number of membranes $M$ becomes clear when plotting $\tau$ of the ladder circuit in Fig. 3F for a stack of 1 to 200 identical POPC bilayers. In case of equal resistance $R_M$ and capacitance $C$, the ladder's largest time constant can be computed analytically[45] as $\tau = R_M C(2 - 2\cos(\pi/M))^{-1}$, which for larger $M$ values can be Taylor approximated as $\tau \approx M^2 \cdot (R_M C/4)$. Concretely, the value of $\tau$ ranges from 29.952 ns for one POPC bilayer up to 505.69 $\mu$s for 200 POPC bilayers. This illustrates the very large increase in oxygen transport time scales when stacking bilayers. The time constant $\tau$ determines the slowest time-scale at which the total oxygen storage in the myelin stack can respond to changes in the extracellular and intracellular oxygen concentrations. However, if, for example, the extracellular oxygen concentration drops at a rate that is faster than the slowest time constant of a 10-bilayer stack, then the bilayers closest to the extracellular side will release a part of their oxygen storage at a rate that matches the extracellular rate of change. This is because the 10-bilayer stack has 9 other time constants that are smaller than $\tau_{\text{ref}}$. To summarize, it is the collective storage of oxygen over all the bilayers in its stack that collectively responds to oxygen concentration changes with a characteristic time equal to $\tau_{\text{max}}$.

### Building an extended RC circuit to model oxygen transport from capillary to myelinated neurons

As oxygen consumption is central to energy production, it is often used to probe neuronal activity. Following neuronal activation, both the BOLD signal in fMRI studies and direct measurements of $P_{O_2}$ using 2PLM ($\leq$100 ms) due to increase in $CMRO_2$ prior to the hyperemic CBF response (>500 ms) [17,18,21–25]. It was recently shown how the initial dip in $P_{O_2}$ is causal to increased local blood flow response (capillary hyperemia), even in the absence of neuronal stimuli[46]. The experimentally observed time scale of $P_{O_2}$ responses is significantly slower than the characteristic time scales of oxygen storage in myelin. This discrepancy arises because experiments typically probe regions in or near capillaries, which are farther from the oxygen-consuming axonal mitochondria than the myelin sheath enveloping the axon.

To bridge the gap between our model and experimental observations, we extended the RC ladder circuit model for a myelin sheath to include the cytosolic solvent inside the axon and the extracellular matrix solvent. This extension allows
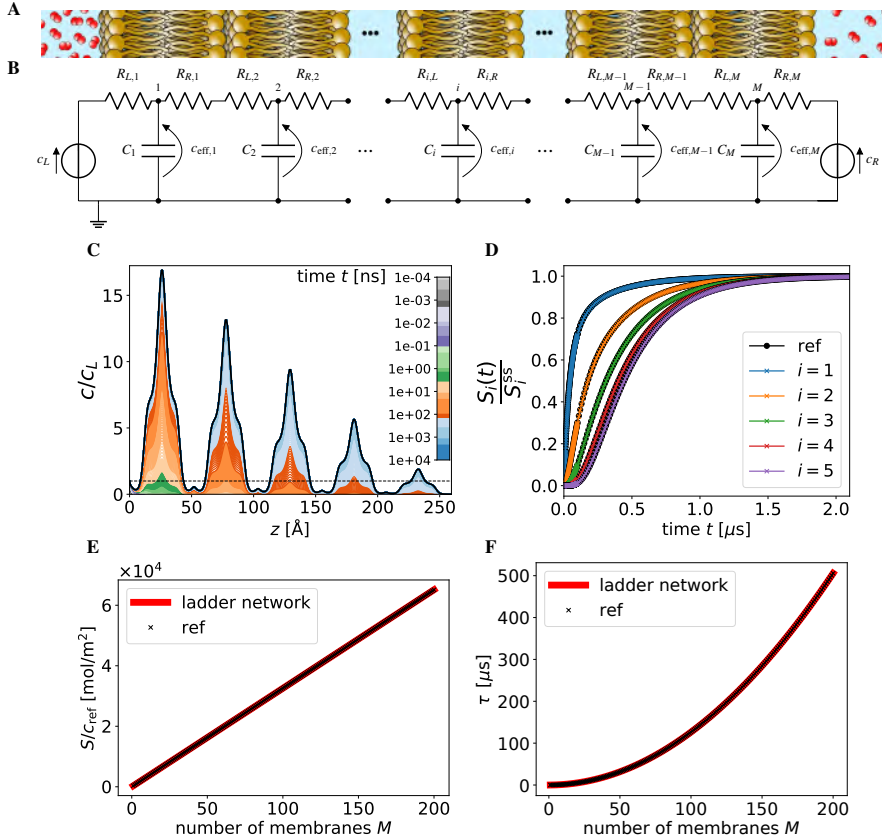
**Figure 3.** Cartoon representation, $O_2$ molecules not to scale. (**A**) and ladder network representation (**B**) of a myelin sheath consisting of $M$ bilayers, separating two solvent regions. Panels **C-D** study an initially oxygen-depleted sheath of $M = 5$ bilayers, where at time $t = 0$ the left solvent oxygen concentration is set to $c_L$. **C**: Oxygen concentration profile $c(t,z)/c_L$ in sheath evolves towards steady-state concentration $c^{ss}(z)$ (thick black line). Dashed horizontal line denotes $c/c_L = 1$. **D**: Normalized oxygen storages $S_i(t)/S_i^{ss}$ (colored lines). Storage $S_i(t)$ in each bilayer $i$ evolves towards its steady-state storage $S_i^{ss}$. **E**: The total oxygen storage (per membrane surface area, when embedded in solvent with oxygen concentration $c_{ref}$) increases linearly with the number of bilayers $M$ it contains. **F**: The largest time constant $\tau$ of oxygen displacement through a myelin sheath depends (approximately) quadratically on the number of bilayers $M$. In panels **E-F**, the red line is the analytical result of the RC ladder network; black markers are values calculated using the discretized Smoluchowski model, serving as reference (ref).

us to study oxygen transport from capillaries to axonal mitochondria, as visualized in Fig. 4. The extension was built by adding (wider) 'solvent' circuit elements on the left and right sides of the ladder circuit (see Methods for parametrization). This extended model was then used to investigate the effect of myelin sheaths on oxygen transport following neuronal activation. We start by describing the two cases that will be investigated with the set-up, after which the choice of specific parameters is discussed.

In the extended set-up, the left concentration $c_L(t)$ now refers to the oxygen concentration near the capillary at $z = 0$ and the right concentration $c_R(t)$ to that in the axonal mito-
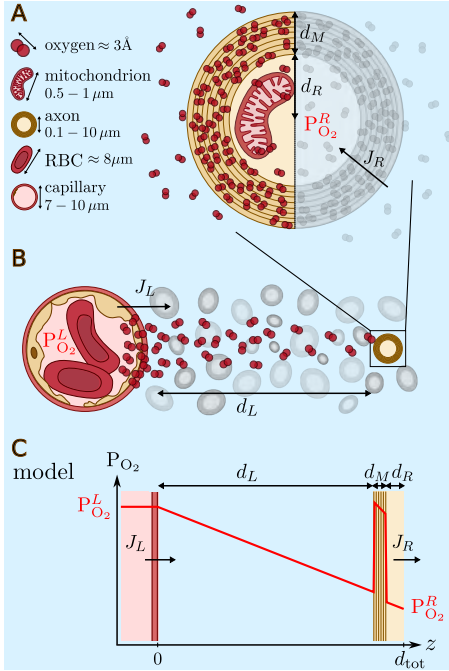
**Figure 4.** Sketches, $O_2$ not to scale. **A**: Sketch of myelinated axon. **B**: Sketch of capillary and myelinated axons. **C**: Model of oxygen transport from a capillary to a myelinated axon. Indicated distances: $d_R$ is the radius of the axon, $d_M$ is the thickness of $M$ bilayers in the myelin sheath, $d_L$ is the distance from the myelin sheath to the capillary. Indicated fluxes: $J_L$ (left) is $O_2$ flux flowing from capillary at $z = 0$, $J_R$ (right) is $O_2$ flux reaching the mitochondria at $z = d_{\text{tot}}$.

chondria at $z = d_{\text{tot}}$ (Fig. 4C). The flux $J_L(t)$ represents the oxygen influx from the capillary, and the flux $J_R(t) = \text{OCR}(t)$ represents the oxygen consumption rate (OCR) at the axonal mitochondria. Generally speaking, the extended set-up will be used to investigate two cases: the case of steady-state oxygen concentrations and the case of the response to a increase in oxygen demand. First, in steady-state, the concentrations $c_L$ and $c_R$ are kept constant in time, which means that the corresponding oxygen partial pressures $P_{O_2}^L$ and $P_{O_2}^R$ are assumed to be constant. The resulting oxygen flux is the steady-state flux $J^{\text{ss}}$, which is the OCR in steady-state ($J_L(t) = J_R(t) = J^{\text{ss}}$, $t < 0$). Second, starting from this steady-state case, it is assumed that action potentials have passed leading to an increase in oxygen demand at $t = 0$. The increased OCR is modeled by a gradually increasing $J^R(t)$ to the mitochondria for $t > 0$

(Fig. 5B). At the capillary side, assuming the blood flow has not changed yet, the oxygen flux is kept at its steady-state value ($J_L(t) = J^{\text{ss}}$, $t > 0$). This case is investigated by imposing constant flux BCs (Fig. 5A). As the OCR increases, the oxygen concentration $c_R(t)$ at the axonal mitochondria will decrease and eventually reach 0. We now define the sustainment time $T$ as the time it takes for $c_R(t)$ to reach 0, which is interpreted as the time that the set-up can sustain the increased oxygen demand.

We varied several parameters in the extended set-up. Most parameters are indeed not uniquely known experimentally at sub-micrometer resolution, but moreover it allows us to explore their effect on the oxygen transport. Variation in geometric parameters and solvent diffusivity gave a total of 30 different configurations for the set-up. The geometric parameters of the extended set-up involve the thickness of the myelin sheath $d_M$, the distance from the myelin sheath to the capillary $d_L$, and the radius of the axon $d_R$. The thickness of the myelin sheath was defined by the number of POPC bilayers $M$ it contains, where $M = 1, 10, 25, 50,$ and $100$ were considered ($d_{M=1} = 5.2\,\text{nm}$). The (inner) axonal radius was set to $d_R = 0.1\,\mu\text{m}, 0.25\,\mu\text{m},$ and $0.5\,\mu\text{m}$, which is representative for most CNS axons[47]. The distance from the myelin sheath to the capillary $d_L$ was adjusted such that all set-ups had the same total length $d_{\text{tot}} = d_L + d_M + d_R \approx 20.53\,\mu\text{m}$. The reference $d_{\text{tot}}$ was chosen for the ($d_{M=100} \approx 0.52\,\mu\text{m}$, $d_R = 0.1\,\mu\text{m}$, $d_L = 20\,\mu\text{m}$) set-up. As $d_L$ is by far the largest contribution to $d_{\text{tot}}$, its value varies closely around $20\,\mu\text{m}$, which was chosen to be close to the mean extravascular distances of $21.4 \pm 0.6\,\mu\text{m}$ reported in human brain capillary networks[48]. The myelin sheath in the set-up was built from the POPC RC circuit element as in Fig. 3. For the parameters to define the 'solvent' RC circuit element for the extracellular and cytosolic solvent, two values for the oxygen diffusion coefficient $D_{\text{solvent}}$ in solvent were considered. The first value was chosen as the diffusion coefficient of oxygen in water $D_w \approx 5.08 \times 10^{-5}\,\text{cm}^2/\text{s}$ as estimated by MD simulations[36]. Secondly, a value of $D_w/2$ was considered to account for overestimation of the diffusivity by the MD simulation[36] and by molecular crowding effects[39,49–51].

Similarly to the variations in the extended set-up parameters, we also varied the boundary conditions allowing us to simulate two interesting cases, i.e. the case of steady-state and the case of increased oxygen demand. For the case of steady-state, the oxygen concentration $c_L(t)$ near the capillary (left) is set to $P_{O_2}^L = 35\,\text{mmHg}$ (conversion $c$ to $P_{O_2}$ is discussed in Methods section). This value is in the range of reported $P_{O_2}$ values in brain tissue[14,52–55]. The oxygen concentration $c_R(t)$ at the axonal mitochondria (right) was varied between $P_{O_2}^R = 1\,\text{mmHg}, 5\,\text{mmHg},$ and $10\,\text{mmHg}$[56,57].

When simulating the case of increased oxygen demand, the strength and time-scale of the increase were varied. An approximately 5-fold OCR increase has been reported in the hippocampal CA3 network during gamma oscillations as compared to absence of spiking, or, equivalently, a 2-fold

increase as compared to spontaneous activity[40]. Likewise, oxygen-consumption rates were found to be about five times higher during seizure activity compared to interictal activity in CA3[58]. Therefore, OCR increase factors $f = 2, 3, 4,$ and 5 were considered in the simulations. Values for the characteristic time $\tau_{OCR}$ were not found. However, $\tau_{OCR}$ is expected to be $\leq 100\,$ms, as it is causal to the initial dip in $P_{O_2}$. Therefore, choices of $\tau_{OCR} = 1\,$ms, $10\,$ms, and $100\,$ms were considered in the simulations. The increase in OCR (OCR $\equiv J^R(t)$) was assumed to follow an exponential shape, starting from OCR$= J^{ss}$ at $t = 0$ reaching OCR $= f \times J^{ss}$ at long times (Fig. 5B).

**Oxygen availability in myelinated axons after neuronal activity**

The results of the simulations using the extended set-up of Fig. 4C are now discussed. The steady-state flux $J^{ss}$ is shown for all configurations in Fig. 5C. Each column corresponds to a value of $P_{O_2}^R$ and each subpanel to a choice for $D_{\text{solvent}}$. The flux ranges from $3.98\,\mu\text{mol/m}^2/\text{s}$ to $11.04\,\mu\text{mol/m}^2/\text{s}$. We describe the effect of the parameters $D_{\text{solvent}}$, $P_{O_2}^R$, $M$, and $d_R$. The flux is impacted mostly by the solvent oxygen diffusivity $D_{\text{solvent}}$. This is understood by $J^{ss} = R_{\text{tot}}^{-1}(c_L - c_R)$, where $R_{\text{tot}}$ is the total resistance of the extended set-up and $c_L - c_R$ is the oxygen concentration difference corresponding to $P_{O_2}^L - P_{O_2}^R$. The solvent makes up most of the total resistance $R_{\text{tot}}$ ($> 90\%$) due to it being the largest phase of the system ($> 97\%$). As the resistance of water to oxygen is proportional to its diffusivity, the flux is approximately twice as large when $D_{\text{solvent}} = D_w$ compared to $D_{\text{solvent}} = D_w/2$. Next, the concentration difference ($P_{O_2}^L - P_{O_2}^R$) impacts $J^{ss}$ linearly. The flux decreases with increasing myelination ($M$), for both solvent diffusivities, as the permeability through the phospholipid bilayers remains lower than the permeability through the solvent phases. The axonal diameter $d_R$ has no impact on $J^{ss}$, as changes in $d_L$ were compensated by changes in $d_R$ (to keep $d_{\text{tot}}$ equivalent).

The sustainment time $T$ after an OCR increase in the axon is given in Fig. 5D for all configurations. Each column corresponds to a choice of $P_{O_2}^R$, each row to a choice for $\tau_{OCR}$, each subplane to a choice of $D_{\text{solvent}}$, and each color a choice of increase factor $f$. Stronger OCR increases (larger $f$) have shorter sustainment times $T$, for all configurations. The axonal radius $d_R$ has negligible impact on $T$, except for the fastest OCR increase at the lowest cytosolic oxygen concentration (top left plot in Fig. 5D). The smaller $D_{\text{solvent}}$ results in longer sustainment, and its impact is more pronounced for the fastest OCR increase. The cytosolic oxygen concentration $P_{O_2}^R$ has a large impact on $T$. This is to be expected, as $P_{O_2}^R$ directly impacts the total oxygen availability at the axonal mitochondria.

Next, we investigate how myelination ($M$) affects the time it takes for the axonal oxygen concentration to drop to zero. On one hand, a thicker myelin sheath has a higher resistance, slowing oxygen permeation from the extracellular matrix to the cytosol and potentially leading to faster depletion at the axonal mitochondria. On the other hand, the capacitance of membranes is higher than that of solvents, where the myelin sheath may function as an oxygen reservoir, counteracting this depletion risk. In our simulations, an increase in myelination prolongs the sustainment $T$ for all configurations. This shows that the buffering effect on $T$ outweighs the negative impact of the higher permeation resistance. Myelin sheaths can thus be viewed as oxygen buffers rather than barriers for oxygen transport.

As expected, none of the set-ups can sustain the OCR increase indefinitely, as the influx $J^L$ is not matching the OCR increase in the simulation. The largest stack ($M = 100$) with the weakest and slowest OCR increase ($f = 2$, $\tau_{OCR} = 100\,$ms) can sustain the oxygen demand for approximately 115 ms for the highest cytosolic oxygen concentration ($P_{O_2}^R = 10\,$mmHg) and the slow oxygen diffusivity $D_{\text{solvent}} = D_w/2$, as visible in Fig. 5D (bottom row, right panel). Given that the vascular response is observed to follow neuronal activity within $\geq 500\,$ms, this results in an approximate 400 ms period of unmet oxygen demand, supporting the functional role of the hyperemic vascular response to restore oxygen levels.

## 3 Discussion

Oxygen transport through a phospholipid bilayer is governed by first-order kinetics, and can be accurately described using an intuitive RC circuit model with membrane resistance ($R_M$) and membrane capacitance ($C$) parameters. The picture of oxygen transport through phospholipid membranes is one of oxygen gradually flowing into the membrane interior, as if an electric capicitor is charged, where the head-group region of the phospholipid membrane creates some resistance to oxygen, as if it were an electric resistor. Phospholipid membranes can thus be loaded and unloaded as if they were an electric RC circuit, and membranes can therefore be regarded as spatially compact oxygen reservoirs.

Myelin, as a stack of phospholipid bilayers, is modeled as a ladder circuit of individual bilayer RC elements. A myelin sheath's resistance to oxygen permeation and its capacitance to store oxygen increase linearly with the number of layers, while the characteristic time constant of its stored oxygen increases quadratically as $\tau \approx M^2(R_M C/4)$. Myelin is primarily composed of lipids (70 % to 85 %, of which 40 % cholesterol, 40 % phospholipids, and 20 % glycolipids) and proteins (15 % to 30 %)[28]. Both the presence of cholesterol and protein are known to reduce membrane permeability to oxygen [29,59]. While our model was constructed from MD data in a pure POPC bilayer, we have shown in the Supplementary Information that the model is applicable to a broad range of bilayers, as it depends solely on the typical shape of the $F$ and $D$ profiles of a non-polar molecule in a bilayer. Therefore, the qualitative results of our model (dependence on myelin thickness) will hold for the more complex bilayers of myelin, while the values of the characteristic time constants would differ. Regardless of the idealized membrane composition, our work shows that myelin has a buffering effect, where the
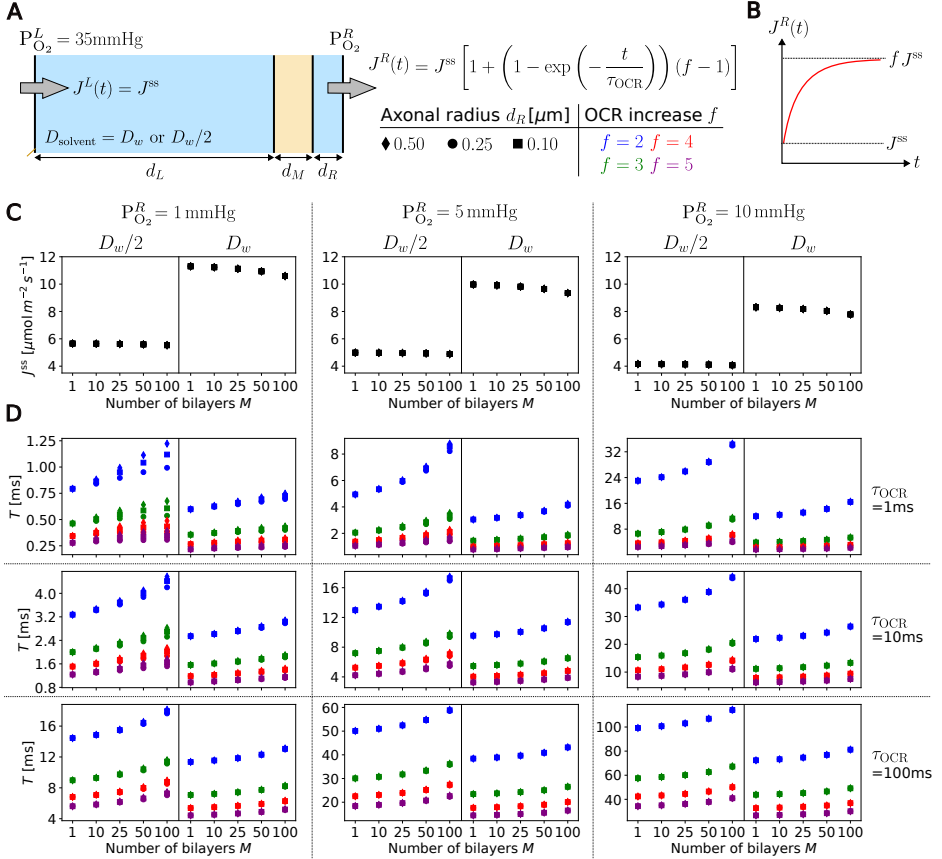
**Figure 5.** Panel **A**: schematic representation of myelin sheath connected to extracellular (left) and cytosolic (right) solvents. $P_{O_2}^L$ and $P_{O_2}^R$ are the oxygen partial pressures near the capillary and axonal mitochondria, respectively. Distances: $d_M$ is the myelin sheath thickness, $d_L \approx 20\,\mu m$ is the distance from the myelin sheath to the capillary, and $d_R$ is the axonal radius. $D_{solvent}$ is the oxygen diffusion coefficient in the solvent phases, where $D_w$ is the diffusion coefficient in water. $J_L$ represents the oxygen influx from the capillary, and is kept to its steady-state value $J^{ss}$ throughout the simulation. Panel **B**: $J_R(t)$ represents the OCR at the axonal mitochondria, and is increased from $J^{ss}$ to $f \times J^{ss}$ using an exponential shape with time constant $\tau_{OCR}$. Panel **C**: the steady-state flux $J^{ss}$ for all configurations as a function of the number of POPC bilayers $M$. The three columns represent different choices for $P_{O_2}^R$ and the subpanels the two choices of $D_{solvent}$. Panel **D**: the sustainment times $T$ for all configurations as a function of the number of POPC bilayers $M$. Columns are different choices for $P_{O_2}^R$. Rows are different choices for $\tau_{OCR}$. Colors denote different OCR increase factors $f$ (2 in blue, 3 in green, 4 in red, and 5 in purple). Marker shapes are according to axonal radius $d_R$ (in µm).

membranes not only act as reservoirs, but also soften sudden changes in oxygen because of the loading and unloading timescales increasing quadratically with the myelin thickness

(from $\approx 30\,$ns for 1 membrane to 506 µs for 200 membranes).

Connecting the myelin membranes to the larger micrometer scale, we modeled neuronal activity by adding large solvent

phases to the myelin RC ladder network. Starting from a steady-state configuration, the OCR at the axon interior was gradually changed to model increased oxygen demand. The steady-state fluxes obtained in our set-ups ranged from 4 to 11 $\mu mol/m^2/s$, and the strongest OCR increase factor $f$ was set to 5. In comparison, experimental measurements of the oxidative capacity of isolated muscle mitochondria were found to be in the range of 2 to 15 $\mu mol/m^2/min$, where the $m^2$ value pertains to mitochondrial surface area[60]. Mitochondria will be sparsely distributed in the cytosol of white matter axons, but considering their large surface to volume ratio[60,61], the OCR values obtained in our set-ups are reasonable.

The oxygen concentration $P_{O_2}^R$ at which the axon operates at inactive (steady-state) conditions has the highest impact on the sustainment time $T$. This concentration will in practice depend on both the amount of capillaries and the distance of those capillaries in the direct vicinity of the axon segment. There is a paradoxical relationship, where sufficient oxygen availability is required to sustain activity, and yet not too much oxygen should be present to avoid oxidative stress. The latter is especially important in the context of neurons, which are prone to the adverse effects of reactive oxygen species [62,63]. Multiple values for $P_{O_2}$ (1, 5, and 10 mmHg) were taken to probe different operational configurations. None of the configurations considered in our model could support the OCR increase until the CBF response is typically observed. On one hand, the insufficiently long myelin buffering timescale could be caused by limitations in the modeling, such as idealized membrane composition, simplified set-up, and inaccurate parameters. On the other hand, the observed oxygen depletion could be supporting the view that the CNS operates at low oxygenation levels and that the 'overshooting' CBF response is essential to restore oxygen homeostasis in the brain[53,64–66].

Besides the idealized membrane composition, another limitation is that the myelin sheath and the solvent phases have been modeled using a one-dimensional description. As such, each point $z$ represents an infinitely extended plane, where the neuronal set-up effectively models oxygen flow from one plane of tissue near a capillary to mitochondria at an 'axonal plane'. Although the model could be extended to a three-dimensional description containing a distribution of cylindrical myelinated axons and capillaries, we argue that the one-dimensional model already captures qualitative aspects of the sustainment time $T$. When the oxygen concentration in the axon interior reaches 0 due to the OCR increase, this does not imply that the oxygen concentration in the myelin and the extracellular matrix is also depleted (i.e. there is still an oxygen flux to the mitochondria, this flux is just lower than the OCR demanded to sustain neuronal activity.) This observation of *local* oxygen depletion near the axonal mitochondria is expected to remain valid in a three-dimensional description.

It is not straightforward to put our results of oxygen buffering into context, given the complex biochemistry in the energy household. In human brains, myelin-rich white matter accounts for 44 % of cortical volume[67]. Disruption of nor-

mal myelin patterns in conditions such as stroke, traumatic brain injury, or multiple sclerosis causes clinical morbidity largely due to damage to central white matter tracts[68]. Psychiatric disorders, including depression and schizophrenia have been associated with white matter defects[69]. Studies in humans and rodents support the existence of lifelong white matter plasticity. Myelin is not only generated early in development, but also during adulthood contributing to learning and memory, and myelination is regulated in response to neuronal activity[70,71]. Training such as extensive piano practicing can induce white matter plasticity in myelinating tracts[72]. In animal experiments, myelin thickness changes have been reported in relation to neuronal stimulation.[73,74]. Following motor learning in the mouse primary motor cortex, addition of newly formed myelin sheaths increases the number of continuous stretches of myelination[71]. Sheath additions may serve to increase conduction velocity, stabilize adaptations in the axonal architecture, or support the increased metabolic demands of neurons[75]. Myelin is crucial for rapid neuronal signal conduction, but also for metabolic support of axons. Myelin in the CNS is generated by oligodendrocytes, and oligodendrocytes synthesize lactate via aerobic glycolysis and transport it into the axon via monocarboxylate transporters, establishing metabolic coupling relevant for both health and disease[10,27,76]. In summary, we have shown here that increased myelination results in longer sustainment of increased oxygen demand, pointing towards a short-time buffer function of myelin. Clearly, the role of myelin as oxygen storage in supporting neuronal activity is a key area for further study.

## 4 Methods

### Smoluchowski diffusion model

Consider a flat membrane with normal along the $z$-axis, that separates two solvent regions with oxygen concentrations $c_L$ ($z < 0$, left) and $c_R$ ($z > h$, right), see Fig. 1A. Averaged over time, the membrane is inhomogeneous along the membrane normal ($z$-axis) and homogeneous in the two directions parallel to the membrane ($x, y$-axes)[36]. If oxygen displacement is assumed to be diffusive, the oxygen concentration $c(z,t)$ is given by the one-dimensional Smoluchowski equation

$$\frac{\partial}{\partial t}c(t,z) = \frac{\partial}{\partial z}\left[D(z)e^{-\beta F(z)}\frac{\partial}{\partial z}\left(e^{\beta F(z)}c(z,t)\right)\right], \quad (1)$$

where $D(z)$ is the position-dependent diffusivity, $F(z)$ is the oxygen free energy profile, and $\beta = (k_B T)^{-1}$ is the inverse temperature with $k_B$ the Boltzmann constant and $T$ temperature. For constant concentrations $c_L$ and $c_R$, the oxygen concentration converges to a steady-state distribution, where the net oxygen flux remains constant. Denoting the concentration difference $c_R - c_L$ by $\Delta c$, the net oxygen flux $J$ can be related to the membrane permeability $P$ by $J = -P\Delta c$. The permeability can be obtained directly from the diffusivity and free energy profiles[38,77]

$$\frac{1}{P} = \int_0^h \frac{dz}{D(z)e^{-\beta(F(z)-F_{ref})}}. \quad (2)$$

where $F_{\text{ref}}$ is the free energy of oxygen in the solvent regions next to the membrane. The inverse of the permeability can be interpreted as a resistance, meaning that the membrane acts as a barrier to permeant diffusion. Moreover, the membrane acts as an oxygen reservoir[34]. The oxygen storage $S(t)$ in a membrane (per unit of cross area in the $xy$-plane) is obtained by integrating the oxygen concentration over the membrane thickness $h$,

$$S(t) = \int_0^h c(z,t)\,\mathrm{d}z \qquad (3)$$

In equilibrium (eq), the solvent contains a reference concentration $c_L = c_R = c_{\text{ref}}$ leading to the equilibrium storage $S^{\text{eq}}$ in the membrane,

$$S^{\text{eq}} = c_{\text{ref}} \int_0^h e^{-\beta(F(z)-F_{\text{ref}})}\,\mathrm{d}z. \qquad (4)$$

In this work, multiple models are built based on the Smoluchowski equation (Eq. 1), which are listed in Fig. 6, and which are discussed in detail below.

## MD simulations and discretized Smoluchowski rate model

In this work, previously published MD simulations of oxygen molecules in a 1-palmitoyl-2-oleoyl-sn-glycero-3-phosphocholine (POPC) bilayer are used[36]. The MD simulations were performed with the CHARMM36 lipid force field[78,79], the modified TIP3P water model, and a custom oxygen model[36]. Simulations were performed at 310 K with a box containing 72 POPC lipids, 5757 waters, and 10 $O_2$ molecules, as visualized in Fig. 1A.

The free energy profile $F$ and diffusivity profile $D$ were extracted from the MD simulation data using Bayesian analysis[80]. Spatial discretization into $N = 100$ bins of width $\Delta z = 0.68$Å along the membrane normal $z$ was used to create a transition matrix $T$ from the MD trajectory data. Then, the discretized $F$ and $D$ profiles with the highest likelihood of generating these observed transitions were constructed using a Monte Carlo procedure. Subsequently, the rate matrix $R$ was constructed from these discretized profiles[80]. The discretized Smoluchowski model gives a set of rate equations $\mathrm{d}c/\mathrm{d}t = R \cdot c$, where $c$ is the $N$-dimensional concentration vector and $R$ is the $N \times N$ tridiagonal rate matrix[81]. The element $c_i$ denotes the oxygen concentration in bin $i$, while the element $R_{i,j}$ denotes the transition rate from bin $j$ to bin $i$. Due to the periodic boundary conditions (PBCs) employed in the MD simulations, the $F$ and $D$ profiles are periodic, resulting in non-zero corner elements $R_{N,1}$ and $R_{1,N}$ in the rate matrix (Model 1, Fig 6). The discretized Smoluchowski model allows us to numerically study the time-dependent concentration profile of oxygen within the POPC bilayer, as the solution is readily available as $c(t) = e^{Rt}c(0)$, where $c(0)$ is the initial concentration vector and $e^{Rt}$ denotes the matrix exponential of $Rt$.

## Adaptation for rate models used in this work

The periodic rate matrix $R$ was used to extract $F$ and $D$ in the Bayesian analysis on the MD trajectories. In this work, $F$, $D$, and/or $R$ have been adapted to model different scenario's: to focus on one bilayer, to construct a stack of bilayers in case of a myelin sheath, and to impose different boundary conditions (BCs).

To focus on the oxygen kinetics within one membrane (Model 2, Fig. 6), the $F$ and $D$ profiles are first truncated to contain only the $N_1 = 76$ membrane bins. A bin in each solvent region is also considered, to be able to impose a constant oxygen concentration in those two water bins. This gives $N_1 + 2 = 78$ bins in total in the $F$ and $D$ profiles, and the corresponding rate matrix[81] $R$ for one bilayer has dimension $78 \times 78$.

A myelin sheath consisting of $M$ bilayers (Model 4, Fig. 6) can be constructed by appending the truncated $F$ and $D$ profiles ($N_1 = 76$ bins) of a single bilayer $M$ times, and then constructing the corresponding rate matrix $R$. A water bin is again included in each solvent region. This results in a free energy profile, diffusion profile, and concentration vector $c$ of dimension $M \cdot N_1 + 2$, and a rate matrix $R$ of dimension $(M \cdot N_1 + 2) \times (M \cdot N_1 + 2)$ for the bilayer stack.

Rather than solving the Smoluchowski equation with PBCs, we aim for studying the time-dependent behavior under oxygen concentrations BCs in the solvent regions. Consider a rate matrix of dimension $N + 2$, where a bin was added to the left (bin 0) and to the right (bin $N + 1$) to impose the solvent concentrations. The $(N+2) \times (N+2)$ rate matrix $R$ with PBC (for one bilayer or a stack of $M$ bilayers) needs to be manipulated to establish constant oxygen concentrations $c_0 = c_L$ and $c_{N+1} = c_R$ in the left and right solvent bins, respectively. This is accomplished by setting the first and last row of the rate matrix $R$ to zero. The zero $R_{N+1,0}$ and $R_{0,N+1}$ elements remove the periodic boundary effects (in case those were imposed), while the zero $R_{0,0}$, $R_{0,1}$, $R_{N+1,N+1}$ and $R_{N+1,N+0}$ elements ensure that $c_0$ and $c_{N+1}$ remain constant. Next, we have shown that the solution of the rate model under such fixed concentration BCs is based on the matrix $R'$, which is the $N \times N$ rate matrix stripped from its first/last rows and columns[34]. With $c'$ the $N$-dimensional vector stripped of its first/last element, and with $U$ an $N$-dimensional vector with zeros (except for $U_1 = c_L$ and $U_N = c_R$), the general solution in case of fixed concentration BCs becomes

$$c'(t) = e^{R't}c'(0) + R'^{-1}\left(e^{R't} - I\right)U. \qquad (5)$$

For one bilayer, this results in Model 2 of Fig. 6 and for a stack of $M$ bilayers, this results in Model 4.

## Kinetic approximation of the rate model for one bilayer

Consider the $N_1 \times N_1$ truncated rate matrix $R'$ obtained by stripping the first and last rows/columns of the rate matrix $R$ of the rate model for one POPC bilayer (Model 2 in Fig. 6). This truncated matrix $R'$ represents the part of $R$ related to the

| model | sketch | rate matrix | # bins | BCs | EQs |
|---|---|---|---|---|---|
| 1 | | $R$ | 100 | PBC | |
| 2 | | $R'$ | $76 + 2$ | conc | 5 |
| 3 | | $\dfrac{1}{\tau_{\mathrm{RC}}}$ | $1 + 2$ | conc | 8, 9 |
| 4 | | $R'$ | $(M \times 76) + 2$ | conc | 5 |
| 5 | | $A, A'$ | $M + 2$ | conc | 10, 11 |
| 6 | | $A^{\mathrm{neur}\prime}$ | $(500 + M + 50) + 2$ | conc | 12 |
| 7 | | $A^{\mathrm{neur}\prime\prime}$ | $(500 + M + 50) + 2$ | flux | 13 |

**Figure 6.** Models used in this work, with their respective rate matrices, number of bins, boundary conditions (BCs), and relevant equations. Model 1: discretized Smoluchowski rate model for one POPC bilayer, built directly from MD simulation data. Model 2: removal of PBCs and solvent bins from Model 1, and introducing concentration BCs at both sides of the bilayer. Model 3: RC circuit model for 1 bilayer with concentration BCs. Model 4: discretized Smoluchowski rate model for a myelin sheath of $M$ bilayers. Model 5: RC ladder network of a myelin sheath of $M$ bilayers. Model 6: the neuronal activity set-up reintroduces the extracellular ('ext') and cytosolic ('cyt') solvent regions using solvent RC elements (500 for extracellular solvent, 50 for cytosolic solvent) surrounding the $M$-stack. Model 7: same as Model 6, but using flux BCs.

POPC bins, whereas the stripped rows/columns are related to the solvent bins. The $N_1 = 76$ time constants $\tau_{R',i}$ associated to $R'$ are given by the negative inverse of the eigenvalues $\lambda_{R',i}$. These time constants are shown Fig. 7A, where the large gap between the largest and second largest time constant suggests that oxygen kinetics within a single POPC bilayer can be well approximated by first-order kinetics.

Next consider the $M \cdot N_1 \times M \cdot N_1$ truncated rate matrix $R'$ obtained by stripping the first and last rows/columns of the rate matrix $R$ from the rate model for a myelin sheath consisting of $M$ bilayers (Model 4 in Fig. 6). There exists a large gap between the largest $M$ time constants and all the other time constants (data not shown). This gap allows the kinetics of oxygen within a stack of $M$ bilayers to be well approximated by a model of $M$-th order kinetics.
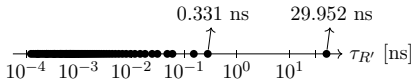


**Figure 7.** Log-scaled plot of the $N_1 = 76$ time constants (black dots) associated to the eigenvalues of $R'$ for a single POPC bilayer.

## Membrane circuit properties

The oxygen capacitance $C$ of POPC is the amount of oxygen stored $S^{\mathrm{eq}}$ (per membrane cross area) in a POPC membrane at equilibrium, per unit of solvent oxygen concentration in the left and right solvent ($c_L = c_R = c_{\mathrm{ref}}$). Using Eq. 4, the capacitance can directly be obtained from the oxygen free energy profile $F(z)$ and the reference free energy $F_{\mathrm{ref}}$ of oxygen in the solvent,

$$C = \int_0^h e^{-\beta(F(z) - F_{\mathrm{ref}})}\, \mathrm{d}z. \tag{6}$$

The resistance $R_M$ of POPC to oxygen permeation is the reciprocal of the membrane permeability $P$, where $P$ is given by Eq. 2.

As the computational models are discretized, the integrals are approximated by sums over the 76 bins associated with POPC

$$C \approx \sum_{i=1}^{76} e^{-\beta(F_i - F_{\mathrm{ref}})} \Delta z = 3.26 \times 10^{-2}\,\mu\mathrm{m}$$
$$R_M = \frac{1}{P} \approx \sum_{i=1}^{76} \frac{\Delta z}{D_i e^{-\beta(F_i - F_{\mathrm{ref}})}} = 3.85\,\mathrm{s/m}, \tag{7}$$

where $F_i$ represents the free energy in bin $i$, while $D_i = (D_{i-1/2} + D_{i+1/2})/2$ is the diffusivity in bin $i$, built from

185

the diffusivity between bins $i-1$, $i$, and $i-1$[36]. The two resistances $R_L$ and $R_R$ are obtained by summing over the first 38 bins and last 38 bins, respectively, in Eq. 7 for $R_M$. Due to the one-dimensional model of the bilayer, the capacitance has a unit length (and not volume), and it connects the oxygen storage $S$ to the solvent concentration $c$ as $S^{eq} = C c_{ref}$ (see main text). For example, a POPC bilayer embedded in water with a 0.5 mg/L $\approx 0.0156$ mol/m$^3$ $\approx 9395$ O$_2$/µm$^3$ oxygen concentration will result in a storage $S^{eq} = 9395$ O$_2$/µm$^3 \times (3.26 \times 10^{-2}$ µm$) \approx 306$ O$_2$/µm$^2$. Thus, at a concentration of 0.5 mg/L, there are approximately 306 oxygen molecules stored within a 1 µm$^2$ slab of a POPC bilayer. Oxygen concentration $c$ can be converted to oxygen partial pressure P$_{O_2}$ as follows. At body temperature and atmospheric pressure, the P$_{O_2}$ is 159 mmHg. Using the oxygen solubility in water ($\sim 7.187$ mg/L[82]), a 10 mmHg P$_{O_2}$ corresponds to $(10/159) \cdot 7.187$ mg/L $= 0.45$ mg/L.

**Time-dependent behavior of RC circuit for one bilayer**
Consider the RC circuit representation (Fig. 2) for a single bilayer (Model 3, Fig. 6). The equation governing the effective potential $c_{eff}(t)$ over the capacitor is given by

$$\frac{dc_{eff}(t)}{dt} = -\frac{1}{\tau_{RC}} c_{eff}(t) + \frac{1}{C} \left( \frac{c_L(t)}{R_L} + \frac{c_R(t)}{R_R} \right), \quad (8)$$

with $\tau_{RC} = C \frac{R_L R_R}{R_L + R_R}$. For a membrane that is initially oxygen-depleted ($c_{eff}(0) = 0$), and for fixed concentration BCs in the solvent regions ($c_L, c_R > 0$, for $t > 0$), the concentration in the membrane evolves as

$$c_{eff}(t) = \frac{R_R c_L + R_L c_R}{R_L + R_R} [1 - \exp(-t/\tau_{RC})] \quad (9)$$

The current through the capacitor denotes the effective oxygen flux $J_{eff}(t)$ through the membrane, and the charge accumulated on the capacitor denotes the oxygen stored within the membrane. Both are obtained via the capacitive relations $J_{eff} = C (dc_{eff}(t)/dt)$ and $S(t) = C c_{eff}(t)$, respectively.

**Solving the ladder circuit for myelin sheath**
A myelin sheath can be modeled as a network of RC circuits (Model 5, Fig. 6). The type of ladder circuit in Fig. 3B is known as a Cauer ladder network. For a sheath of $M$ bilayers, the $M$-dimensional vector $c_{eff}(t)$ is used to collect the effective concentrations $c_{eff}(t) = [c_{eff,1}, ..., c_{eff,M}]^T$. Kirchhoff's current law can be used in each loop of the RC ladder network to construct a set of $M$ first order, linear, and time-invariant differential equations for $c_{eff}(t)$, written in matrix notation as

$$\frac{dc_{eff}}{dt} = A c_{eff}(t) + C^{-1} B U(t) \quad (10)$$

with $A = C^{-1} T$. Here, the $2 \times 1$ matrix $U(t) = [c_L(t), c_R(t)]^T$ denotes the input vector of left and right solvent concentrations. The $M \times M$ diagonal matrix $C$ contains the membrane capacitances, $C = diag([C_1, ..., C_M])$. The $M$x$M$ coefficient matrix $T$ is built from all resistances $R_{L,i}$ and $R_{R,i}$

($i \in [1, ..., M]$) denoting the resistance of the left and right leaflet of bilayer $i$, respectively. It is tridiagonal and symmetric with off-diagonal elements $T_{i,i+1} = T_{i+1,i} = (R_{L,i} + R_{R,i+1})^{-1}$, $i = 1 ... M - 1$. The diagonal elements are $T_{i,i} = -T_{i-1,i} - T_{i+1,i}$, $i = 2 ... M - 2$, and the two corner elements are $T_{1,1} = -1/R_{L,1} - T_{1,2}$ and $T_{MM} = -1/R_{R,M} - T_{M-1,M}$. The $M \times 2$ matrix $B$ contains zeros except for the corner elements, $B_{1,1} = 1/R_{L,1}$ and $B_{M,2} = 1/R_{R,M}$, referring to the permeability for the outermost membrane leaflets. The time behavior of the ladder network is determined by the matrix $A = C^{-1} T$. The $M$ time constants follow from the eigenvalues of $A$. An analytical expression for the eigenvalues exists for the case of identical resistances and capacitances throughout the network[45], which was used to plot the time constants for a sheath consisting of $M$ POPC bilayers (Fig. 3F).

Similar to the RC circuit for one membrane, the storage $S_i$ in each membrane $i$ in the stack and the net flux can be determined from each $c_{eff,i}$ using the capacitive laws $S_i = C_i c_{eff,i}$ and $J_{eff,i} = C_i (dc_{eff,i}/dt)$. Under fixed concentration BCs for $c_L$ and $c_R$ ($U(t) = U$), the general solution for the effective concentrations is given by

$$c_{eff}(t) = e^{At} c_{eff}(0) + A^{-1} (e^{At} - 1) C^{-1} B U \quad (11)$$

An alternative way to solve for concentration BCs in the first and last bins can be implemented by a common technique to 'absorb' the boundary conditions into the rate matrix $A$. Two solvent bins are added to the myelin system, i.e. bin 0 and bin $M + 1$. This results in the $M + 2$ dimensional vector $c'_{eff}$, where $c'_{eff,0} = c_L$ and $c'_{eff,M+1} = c_R$, and the other elements are copied from $c_{eff,i}$. The corresponding $(M + 2) \times (M + 2)$ dimensional rate matrix $A'$ is a copy of matrix $A$, except that it has additional first and last rows/columns. The rows $A'_{0,*}$ and $A'_{M+1,*}$ contain zeros to impose that $c'_{eff,0}$ and $c'_{eff,M+1}$ are constant. The first and last column of $A'$ are also zeros, except for elements $A'_{1,0} = B_{1,1}/C_{1,1}$ and $A'_{M,M+1} = B_{M,2}/C_{M,M}$, thus dictating the flow of concentration from the solvent bins into the network. Using these notations, Eq. 10 with constant concentration BCs is equivalent to solving

$$\frac{dc'_{eff}(t)}{dt} = A' c'_{eff}(t) \quad (12)$$

**Building an RC ladder model for neuronal activity**
In order to model neuronal activity, a larger model is built that comprises a myelin sheath with $M$ bilayers, extracellular solvent, and cytosol solvent (see Fig. 4C). As shown in Model 6 and 7 of Fig. 6, the stack of membranes is modeled by an RC ladder, and solvent is modeled in a similar fashion by a ladder of 'solvent' RC circuits. A 'solvent' RC circuit element has resistances $R_L = R_R = h/(2D_{solvent})$ and capacitance $C = h$, where $D_{solvent}$ is the oxygen diffusivity in the solvent and $h$ is the thickness that one solvent RC element represents. Concretely, the extracellular solvent ($d_L \approx 20$ µm) was represented by a ladder of 500 solvent RC circuits and the cytosol ($d_R = 0.1, 0.25,$ or $0.5$ µm) by 50 solvent RC circuits.

**13/17**

186

For example, for $d_R = 0.5\,\mu m$, the 50 cytosolic solvent RC circuits each represent a solvent patch of thickness $h = 10\,nm$. For $d_L = 20\,\mu m$, the 500 extracellular solvent patches have thickness $h = 40\,nm$. The total resistance $R_{tot}$ of the circuit is the sum of all circuit element resistances.

The rate equations of the extended set-up are similar to Eq. 10, where matrices $C^{neur}$, $T^{neur}$ and $A^{neur}$ have dimension $(M + 550) \times (M + 550)$, and $B^{neur}$ has dimension $(M + 550) \times 2$. Alternatively, two extra solvent bins are added to the ladder with index 0 and $M + 550 + 1$, which aids to impose various boundary conditions.

Under concentration BCs (Model 6, Fig. 6), the rate equations are similar to Eq. 12, where the adapted matrix $A^{neur\prime}$ (constructed similarly as $A'$) has dimension $(M + 550 + 2) \times (M + 550 + 2)$. In the first case of simulations of the extended set-up, the steady-state oxygen concentration profile $c_{eff}^{ss,neur\prime}$ is determined under constant concentration BCs, i.e. $c_L$ and $c_R$ are kept constant. It can be extracted from the sum of the two eigenvectors associated to the $\lambda = 0$ eigenvalues of $A^{neur\prime}$ (corresponding to zero oxygen concentration at the left or at the right). The steady-state flux $J^{ss}$ follows from $J^{ss} = (c_L - c_R)/R_{tot}$.

To impose flux BCs (Model 7, Fig. 6), the first and last rows of $A^{neur\prime}$, which are zero rows, need to be modified further, resulting in rate matrix $A^{neur\prime\prime}$. Moreover, a new input term needs to be introduced to impose the flux BCs (similar to the term $C^{-1}BU$ in Eq. 10 for concentration BCs). Adding such an input term $EJ(t)$, the rate equations for the $M + 550 + 2$ dimensional concentration vector $c_{eff}^{neur\prime\prime}(t)$ become

$$\frac{dc_{eff}^{neur\prime\prime}(t)}{dt} = A^{neur\prime\prime} c_{eff}^{neur\prime\prime}(t) + EJ(t) \qquad (13)$$

Here, the changes in the rate matrix are: $A_{0,0}^{neur\prime\prime} = -A_{1,0}^{neur\prime}$ and $A_{0,1}^{neur\prime\prime} = A_{1,0}^{neur\prime}$. The $(M + 550 + 2) \times 2$ matrix $E$ contains zeros except for the corner elements that are based on the capacitances of the two extra solvent bins, i.e. $E_{0,0} = 1/C_0$ and $E_{M+550+1,1} = 1/C_{M+550+1}$. The input vector $J(t) = [J^L(t), J^R(t)]^T$ contains the fluxes that are imposed at the boundaries.

In the second case of simulations of the extended set-up, the flux near the capillary $J^L(t) = J^{ss}$ is kept constant. Meanwhile, the flux at the axonal mitochondria $J^R(t)$ is increased (expression in Fig. 5B) from $J^R(0) = J^{ss}$ towards a plateau value of $J^R(\infty) = f \times J^{ss}$. The initial concentration vector is the steady-state $c_{eff}^{ss,neur\prime}$. The time evolution of the concentration profile was modeled by feeding the system of rate equations into the Python package 'control'[83].

## References

1. Hyder, F., Rothman, D. L. & Bennett, M. R. Cortical energy demands of signaling and nonsignaling components in brain are conserved across mammalian species and activity levels. *Proc. Natl. Acad. Sci.* **110**, 3549–3554 (2013).

2. Watts, M. E., Pocock, R. & Claudianos, C. Brain energy and oxygen metabolism: emerging role in normal function and disease. *Front. molecular neuroscience* **11**, 216 (2018).

3. Murphy, M. P. How mitochondria produce reactive oxygen species. *Biochem. journal* **417**, 1–13 (2009).

4. Ravera, S. *et al.* Evidence for aerobic atp synthesis in isolated myelin vesicles. *The international journal biochemistry cell biology* **41**, 1581–1591 (2009).

5. Ravera, S., Panfoli, I., Aluigi, M. G., Calzia, D. & Morelli, A. Characterization of myelin sheath f o f 1-atp synthase and its regulation by if 1. *Cell biochemistry biophysics* **59**, 63–70 (2011).

6. Ravera, S., Morelli, A. M. & Panfoli, I. Myelination increases chemical energy support to the axon without modifying the basic physicochemical mechanism of nerve conduction. *Neurochem. Int.* **141**, 104883 (2020).

7. Ravera, S., Bartolucci, M., Calzia, D., Morelli, A. M. & Panfoli, I. Efficient extra-mitochondrial aerobic atp synthesis in neuronal membrane systems. *J. Neurosci. Res.* **99**, 2250–2260 (2021).

8. Morelli, A. M., Chiantore, M., Ravera, S., Scholkmann, F. & Panfoli, I. Myelin sheath and cyanobacterial thylakoids as concentric multilamellar structures with similar bioenergetic properties. *Open Biol.* **11**, 210177 (2021).

9. Morelli, A. M. & Scholkmann, F. Should the standard model of cellular energy metabolism be reconsidered? possible coupling between the pentose phosphate pathway, glycolysis and extra-mitochondrial oxidative phosphorylation. *Biochimie* **221**, 99–109 (2024).

10. Lee, Y. *et al.* Oligodendroglia metabolically support axons and contribute to neurodegeneration. *Nature* **487**, 443–448 (2012).

11. Pooya, S. *et al.* The tumour suppressor lkb1 regulates myelination through mitochondrial metabolism. *Nat. communications* **5**, 4993 (2014).

12. Domenech-Estévez, E. *et al.* Distribution of monocarboxylate transporters in the peripheral nervous system suggests putative roles in lactate shuttling and myelination. *J. Neurosci.* **35**, 4151–4156 (2015).

13. Kim, S. *et al.* Schwann cell o-glcnac glycosylation is required for myelin maintenance and axon integrity. *J. Neurosci.* **36**, 9633–9646 (2016).

14. Erecińska, M. & Silver, I. A. Tissue oxygen tension and brain sensitivity to hypoxia. *Respir. Physiol.* **128**, 263–276 (2001).

15. Krab, K., Kempe, H. & Wikström, M. Explaining the enigmatic km for oxygen in cytochrome c oxidase: A kinetic model. *Biochimica et Biophys. Acta (BBA) - Bioenerg.* **1807**, 348–358 (2011).

16. Gnaiger, E., Lassnig, B., Kuznetsov, A., Rieger, G. & Margreiter, R. Mitochondrial Oxygen Affinity, Respiratory Flux Control And Excess Capacity Of Cytochrome c Oxidase. *J. Exp. Biol.* **201**, 1129–1139 (1998).

17. Buxton, R. B., Uludağ, K., Dubowitz, D. J. & Liu, T. T. Modeling the hemodynamic response to brain activation. *NeuroImage* **23**, S220–S233 (2004). Mathematics in Brain Imaging.

18. Buxton, R. Interpreting oxygenation-based neuroimaging signals: the importance and the challenge of understanding brain oxygen metabolism. *Front. neuroenergetics* **2**, 1648 (2010).

19. Wang, T. *et al.* Hemodynamic response function in brain white matter in a resting state. *Cereb. cortex communications* **1**, tgaa056 (2020).

20. Schilling, K. G. *et al.* Anomalous and heterogeneous characteristics of the bold hemodynamic response function in white matter. *Cereb. Cortex Commun.* **3**, tgac035 (2022).

21. Hillman, E. M. Coupling mechanism and significance of the bold signal: a status report. *Annu. review neuroscience* **37**, 161–181 (2014).

22. Devor, A. *et al.* Coupling of total hemoglobin concentration, oxygenation, and neural activity in rat somatosensory cortex. *Neuron* **39**, 353–359 (2003).

23. Parpaleix, A., Houssen, Y. G. & Charpak, S. Imaging local neuronal activity by monitoring po2 transients in capillaries. *Nat. medicine* **19**, 241–246 (2013).

24. Sakadžić, S. *et al.* Two-photon microscopy measurement of cerebral metabolic rate of oxygen using periarteriolar oxygen concentration gradients. *Neurophotonics* **3**, 045005–045005 (2016).

25. Lecoq, J. *et al.* Odor-evoked oxygen consumption by action potential and synaptic transmission in the olfactory bulb. *J. Neurosci.* **29**, 1424–1433 (2009).

26. Nave, K.-A. & Werner, H. B. Myelination of the nervous system: mechanisms and functions. *Annu. review cell developmental biology* **30**, 503–533 (2014).

27. Simons, M. & Nave, K.-A. Oligodendrocytes: myelination and axonal support. *Cold Spring Harb. perspectives biology* **8**, a020479 (2016).

28. Kister, A. & Kister, I. Overview of myelin, major myelin lipids, and myelin-associated proteins. *Front. Chem.* **10**, 1041961 (2023).

29. Pias, S. C. How does oxygen diffuse from capillaries to tissue mitochondria? barriers and pathways. *The J. Physiol.* **599**, 1769–1782 (2021).

30. Moller, M. *et al.* Direct measurement of nitric oxide and oxygen partitioning into liposomes and low density lipoprotein. *J. Biol. Chem.* **280**, 8850–8854 (2005).

31. Möller, M. N. *et al.* Solubility and diffusion of oxygen in phospholipid membranes. *Biochimica et biophysica acta (bba)-biomembranes* **1858**, 2923–2930 (2016).

32. Möller, M. N. & Denicola, A. Diffusion of nitric oxide and oxygen in lipoproteins and membranes studied by pyrene fluorescence quenching. *Free. Radic. Biol. Medicine* **128**, 137–143 (2018).

33. Al-Abdul-Wahid, M. S. *et al.* A combined nmr and molecular dynamics study of the transmembrane solubility and diffusion rate profile of dioxygen in lipid bilayers. *Biochemistry* **45**, 10719–10728 (2006).

34. Vervust, W. & Ghysels, A. Oxygen storage in stacked phospholipid membranes under an oxygen gradient as a model for myelin sheaths. In *Oxygen Transport to Tissue XLIII*, 301–307 (Springer, 2022).

35. Missner, A. & Pohl, P. 110 years of the meyer–overton rule: predicting membrane permeability of gases and other small compounds. *ChemPhysChem* **10**, 1405–1414 (2009).

36. Ghysels, A., Venable, R. M., Pastor, R. W. & Hummer, G. Position-dependent diffusion tensors in anisotropic media from simulation: oxygen transport in and through membranes. *J. chemical theory computation* **13**, 2962–2976 (2017).

37. Riccardi, E., Krämer, A., van Erp, T. S. & Ghysels, A. Permeation rates of oxygen through a lipid bilayer using replica exchange transition interface sampling. *The J. Phys. Chem. B* **125**, 193–201 (2020).

38. De Vos, O. *et al.* Membrane permeability: Characteristic times and lengths for oxygen and a simulation-based test of the inhomogeneous solubility-diffusion model. *J. chemical theory computation* **14**, 3811–3824 (2018).

39. Place, T. L., Domann, F. E. & Case, A. J. Limitations of oxygen delivery to cells in culture: An underappreciated problem in basic and translational research. *Free. Radic. Biol. Medicine* **113**, 311–322 (2017).

40. Huchzermeyer, C., Berndt, N., Holzhütter, H.-G. & Kann, O. Oxygen consumption rates during three different neuronal activity states in the hippocampal ca3 network. *J. Cereb. Blood Flow & Metab.* **33**, 263–271 (2013).

41. Hodgkin, A. L. & Huxley, A. F. A quantitative description of membrane current and its application to conduction and excitation in nerve. *The J. physiology* **117**, 500 (1952).

42. Kang, Y. S., Hong, J.-M., Jang, J. & Kim, U. Y. Analysis of facilitated transport in solid membranes with fixed site carriers 1. single rc circuit model. *J. membrane science* **109**, 149–157 (1996).

43. Aït-Idir, W. *et al.* Oxygen diffusion impedance in proton exchange membrane fuel cells–insights into electrochemical impedance spectra and equivalent electrical circuit modeling. *Electrochimica Acta* **472**, 143430 (2023).

44. Galpayage Dona, K. N. U., Liu, J., Qiang, Y., Du, E. & Lau, A. Electrical equivalent circuit model of sickle cell. In *ASME International Mechanical Engineering Congress and Exposition*, vol. 58455, V010T13A029 (American Society of Mechanical Engineers, 2017).

45. Yueh, W.-C. Eigenvalues of several tridiagonal matrices. *Appl. Math. E-Notes [electronic only]* **5**, 66–74 (2005).

46. Wei, H. S. *et al.* Erythrocytes are oxygen-sensing regulators of the cerebral microcirculation. *Neuron* **91**, 851–862 (2016).

47. Liewald, D., Miller, R., Logothetis, N., Wagner, H.-J. & Schüz, A. Distribution of axon diameters in cortical white matter: an electron-microscopic study on three human brains and a macaque. *Biol. cybernetics* **108**, 541–557 (2014).

48. Smith, A. F. *et al.* Brain capillary networks across species: a few simple organizational requirements are sufficient to reproduce both structure and function. *Front. physiology* **10**, 233 (2019).

49. Angles, G. & Pias, S. C. Discerning membrane steady-state oxygen flux by monte carlo markov chain modeling. In *Oxygen Transport to Tissue XLII*, 137–142 (Springer, 2021).

50. Ganfield, R., Nair, P. & Whalen, W. Mass transfer, storage, and utilization of o2 in cat cerebral cortex. *Am. J. Physiol. Content* **219**, 814–821 (1970).

51. Homer, L. D., Shelton, J. B. & Williams, T. J. Diffusion of oxygen in slices of rat brain. *Am. J. Physiol. Integr. Comp. Physiol.* **244**, R15–R22 (1983).

52. Tsai, A. G., Cabrales, P. & Intaglietta, M. The physics of oxygen delivery: facts and controversies. *Antioxidants & Redox Signal.* **12**, 683–691 (2010).

53. Leithner, C. & Royl, G. The oxygen paradox of neurovascular coupling. *J. Cereb. Blood Flow & Metab.* **34**, 19–29 (2014).

54. Wilson, D. F. Quantifying the role of oxygen pressure in tissue function. *Am. J. Physiol. Circ. Physiol.* **294**, H11–H13 (2008).

55. Wilson, D. F., Harrison, D. K. & Vinogradov, S. A. Oxygen, ph, and mitochondrial oxidative phosphorylation. *J. Appl. Physiol.* **113**, 1838–1845 (2012).

56. Rasmussen, P. *et al.* Capillary-oxygenation-level-dependent near-infrared spectrometry in frontal lobe of humans. *J. Cereb. Blood Flow & Metab.* **27**, 1082–1093 (2007).

57. Rasmussen, P. *et al.* Cerebral oxygenation is reduced during hyperthermic exercise in humans. *Acta physiologica* **199**, 63–70 (2010).

58. Schoknecht, K. *et al.* Event-associated oxygen consumption rate increases ca. five-fold when interictal activity transforms into seizure-like events in vitro. *Int. J. Mol. Sci.* **18**, 1925 (2017).

59. Dotson, R. J. & Pias, S. C. Reduced oxygen permeability upon protein incorporation within phospholipid bilayers. *Oxyg. Transp. to Tissue XL* 405–411 (2018).

60. Schwerzmann, K., Hoppeler, H., Kayar, S. R. & Weibel, E. R. Oxidative capacity of muscle and mitochondria: correlation of physiological, biochemical, and morphometric characteristics. *Proc. Natl. Acad. Sci.* **86**, 1583–1587 (1989).

61. Schmiedl, A. *et al.* The surface to volume ratio of mitochondria, a suitable parameter for evaluating mitochondrial swelling: Correlations during the course of myocardial global ischaemia. *Virchows Arch. A* **416**, 305–315 (1990).

62. Ndubuizu, O. & LaManna, J. C. Brain tissue oxygen concentration measurements. *Antioxidants & redox signaling* **9**, 1207–1220 (2007).

63. Popa-Wagner, A., Mitran, S., Sivanesan, S., Chang, E. & Buga, A.-M. Ros and brain diseases: the good, the bad, and the ugly. *Oxidative medicine cellular longevity* **2013**, 963520 (2013).

64. Devor, A. *et al.* "overshoot" of o2 is required to maintain baseline tissue oxygenation at locations distal to blood vessels. *J. Neurosci.* **31**, 13676–13681 (2011).

65. Beinlich, F. R. *et al.* Oxygen imaging of hypoxic pockets in the mouse cerebral cortex. *Science* **383**, 1471–1478 (2024).

66. Herculano-Houzel, S. & Rothman, D. L. From a demand-based to a supply-limited framework of brain metabolism. *Front. Integr. Neurosci.* **16**, 818685 (2022).

67. Laughlin, S. B. & Sejnowski, T. J. Communication in neuronal networks. *Science* **301**, 1870–1874 (2003).

68. Micu, I. *et al.* The molecular physiology of the axo-myelinic synapse. *Exp. neurology* **276**, 41–50 (2016).

69. Fields, R. D. White matter in learning, cognition and psychiatric disorders. *Trends neurosciences* **31**, 361–370 (2008).

70. Demerens, C. *et al.* Induction of myelination in the central nervous system by electrical activity. *Proc. Natl. Acad. Sci.* **93**, 9887–9892 (1996).

71. Bacmeister, C. M. *et al.* Motor learning drives dynamic patterns of intermittent myelination on learning-activated axons. *Nat. neuroscience* **25**, 1300–1313 (2022).

72. Bengtsson, S. L. *et al.* Extensive piano practicing has regionally specific effects on white matter development. *Nat. neuroscience* **8**, 1148–1150 (2005).

73. Gibson, E. M. *et al.* Neuronal activity promotes oligodendrogenesis and adaptive myelination in the mammalian brain. *Science* **344**, 1252304 (2014).

74. Mitew, S. *et al.* Pharmacogenetic stimulation of neuronal activity increases myelination in an axon-specific manner. *Nat. communications* **9**, 306 (2018).

189

75. Xin, W. & Chan, J. R. Motor learning revamps the myelin landscape. *Nat. neuroscience* **25**, 1251–1252 (2022).

76. Fünfschilling, U. *et al.* Glycolytic oligodendrocytes maintain myelin and long-term axonal integrity. *Nature* **485**, 517–521 (2012).

77. Marrink, S. J. & Berendsen, H. J. Permeation process of small molecules across lipid membranes studied by molecular dynamics simulations. *The J. Phys. Chem.* **100**, 16729–16738 (1996).

78. Brooks, B. R. *et al.* Charmm: the biomolecular simulation program. *J. computational chemistry* **30**, 1545–1614 (2009).

79. Klauda, J. B. *et al.* Update of the charmm all-atom additive force field for lipids: validation on six lipid types. *The journal physical chemistry B* **114**, 7830–7843 (2010).

80. Hummer, G. Position-dependent diffusion coefficients and free energies from Bayesian analysis of equilibrium and replica molecular dynamics simulations. *New J. Phys.* **7**, 34 (2005).

81. Bicout, D. J. & Szabo, A. Electron transfer reaction dynamics in non-Debye solvents. *J. Chem. Phys.* **109**, 2325–2338 (1998).

82. Battino, R. Solubility data series. *Oxyg. ozone* **7**, 412 (1981).

83. Fuller, S. *et al.* The python control systems library (python-control). In *60th IEEE Conference on Decision and Control (CDC)*, 4875–4881 (IEEE, 2021).

## Acknowledgements

## Author contributions statement

A.G. conceived the concept. W.V. and A.G. developed the theory, implemented the equations, analyzed the results, and wrote the paper. W.V. implemented and ran the computational neuronal activity simulations. K.W. and L.L. provided physiological context in the discussion. All authors reviewed the manuscript.

## Supplementary information

Supplementary Information is available. The scripts used in this work are available on GitHub.

## Additional information

**Competing interests:** The authors declare no competing interests.

# Supporting information:
# Myelin sheaths can act as compact temporary oxygen storage units as modeled by an electrical RC circuit model

**Wouter Vervust[1], Katja Witschas[2], Luc Leybaert[2], and An Ghysels[1,*]**

[1]IBiTech - BioMMeda group, Ghent University, Belgium
[2]Physiology Group, Department of Basic and Applied Medical Sciences, Ghent University, Belgium
[*]an.ghysels@ugent.be

## Contents

## S1 Smoluchowski to RC

A general formula for the time-dependent permeant storage $S(t)$ of non-polar permeants inside a lipid bilayer is now derived. The resulting model is restricted to non-polar permeants, where the free energy profile $F(z)$ contains barriers in the phospholipid headgroup regions ($z \in \{0, L\}$ and $z \in \{h - L, h\}$, Supp. Fig. 1) and a well in the hydrophobic core ($z \in \{L, h - L\}$). The resulting model depends solely on this shape of the free energy profile (barrier-well-barrier), and does not require the $F$ and $D$ profiles to be symmetric w.r.t. the bilayer midplane, nor does it require the solvent regions to have the same free energy. In other words, the model holds for a asymmetric distributions of phospholipids among the bilayer leaflets, and for different solvents at the left- and right sides of the bilayer.

### S1.1 Assumptions

Two assumptions are introduced. The first pertains to a quasi-equilibrium approximation of the permeant concentration within the hydrophobic core, and the second pertains to a quasi-steady-state approximation of oxygen flux through the polar phospholipid head-groups.

The first assumption arises from the steep energy well of the hydrophobic core, where permeant molecules will quickly relax after entering the membrane interior. As such, the permeant concentration of the membrane interior can be approximated to be Boltzmann-distributed: $c(z,t) \propto \exp(-F(z)/(k_B T))$. While this holds true only for the membrane interior, we can extend this assumption to the entire membrane. To see this, we notice that $F(z)$ is largest in the phospholipid headgroup regions, where the corresponding oxygen concentration will be lowest. As an example, the left and right phospholipid head-group regions of POPC both contribute no more than a few percentages of the total equilibrium oxygen storage $S$, as shown in Supp. Fig. 2A. Thus, as

the interior part of the head-group region still contributes to the rising edge of the free energy well, and the exterior part carries negligible permeant concentration, we can extend the quasi-equilibrium approximation to the entire membrane.

The quasi-equilibrium approximation states that the concentration is Boltzmann distributed at all times within the membrane, even when the concentration and the boundary conditions are not implying equilibrium. With $S(t)$ the total amount of permeants in the membrane at time $t$,

$$S(t) = \int_0^h c(z,t)\mathrm{d}z, \tag{S1}$$

this approximation can be written as

$$c(z,t) = S(t)\frac{e^{-\beta\Delta F(z)}}{Z}, \tag{S2}$$

where $\Delta F(z) = F(z) - F_{\text{ref}}$ denotes the free energy difference w.r.t. a chosen reference free energy $F_{\text{ref}}$, and where $Z$ is the partition function ($= C$ the membrane capacitance)

$$Z = \int_0^h e^{-\beta\Delta F(z)}\mathrm{d}z, \tag{S3}$$

This concentration profile will indeed both give the correct total storage $S(t)$ and relative concentration within the membrane.

Using this assumption, we can denote the membrane as a single state $M$, defined by an average membrane concentration $c_M(t)$ and average Boltzmann factor $e^{-\beta\Delta F_M}$, where
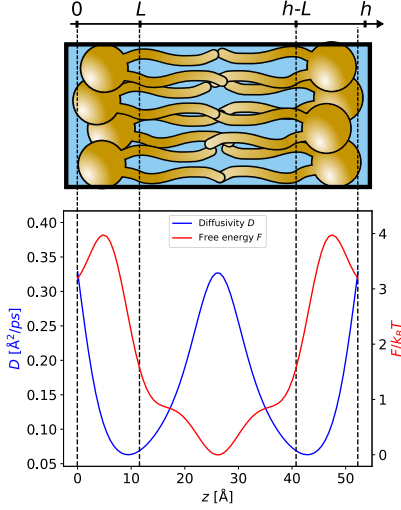
$$\begin{aligned} e^{-\beta\Delta F_M} &= \frac{Z}{h}, \\ c_M(t) &= \frac{S(t)}{h}. \end{aligned} \tag{S4}$$

Note that, using Eq. S2, the above equation results in the following relation, which will be used in the second approximation

$$c_M(t)\, e^{\beta\Delta F_M} = c(z,t)\, e^{\beta\Delta F(z)}. \tag{S5}$$

For the second assumption, we notice that, apart from having a very small contribution to the total permeant storage, the permeant concentration profile within the head-group regions converges significantly faster than the membrane interior. To see this, we look at the scenario of charging an
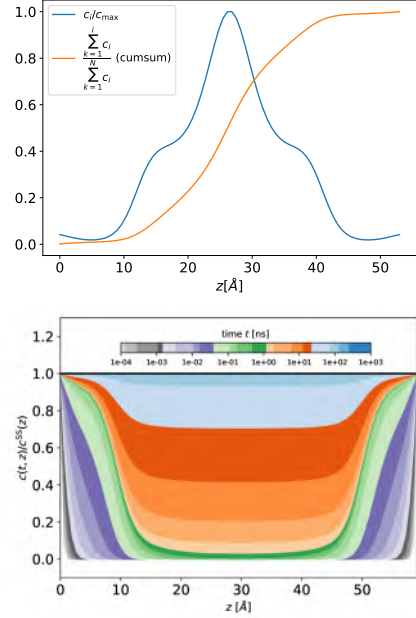
**Supplementary Figure 1.** Free energy (red) and
diffusivity (blue) of oxygen in POPC. The region
$z \in \{L, h-L\}$ denotes the hydrophobic interior of the bilayer,
wheras the regions $z \in \{0, L\}$ and $z \in \{L-h, h\}$ denote the
phospholipid headgroups regions. The specific value of $L$ is
hard to define, as the transition from interior to headgroup
region is rather smooth in nature. $L$ lies somewhere in
between $z = 10\text{Å}$ and $z = 13\text{Å}$.

initially oxygen-depleted POPC bilayer ($c(z, t = 0) = 0$) from
both the left and right ($c_L = c_R = c_{\text{ref}}$, for $t \geq 0$). Supp.
Fig. 2B shows the evolution of $c(z, t)$ through time, where
the concentration profile is divided by the equilibrium pro-
file $c(z, t \to \infty)$ to which it converges. It is seen that the
concentration within the head-group regions quickly rises to
its equilibrium profile. In other words, after a short initial-
ization time $t_{\text{init}}$, the head-group regions absorb a negligible
amount of oxygen. This results in the second assumption,
where the flux $J(z, t)$ throughout the head-group regions is
defined to be constant for $t > t_{\text{init}}$ (e.g. for the left head-group
region $J(z_a, t) \equiv J(z_b, t), \forall z_a, z_b \in [0, L]$). If it were not con-
stant, then the amount of oxygen entering the head-group
region would not equal the amount of oxygen exiting that
region, resulting in a net absorbance (or generation). The flux
through a head-group region can thus be approximated to be
at a quasi-steady-state.

The assumptions are now summarized below. These as-
sumptions normally hold only after a short time $t_{\text{init}}$, but are
approximated to hold for all $t$.

1. The membrane is considered to be at quasi-equilibrium



**Supplementary Figure 2. Top**: equilibrium concentration
of oxygen within POPC, divided by the maximum
concentration at the membrane midplane (blue line).
Normalized cumulative sum of the oxygen concentration
(orange line). Notice how the contribution of the
phospholipid headgroups (starting at $z = 0$ and ending
somewhere in between $z = 10$ and $z = 13$) is only a few
percentages. **Bottom**: Propagation of the oxygen
concentration profile for a POPC bilayer that is initially
oxygen depleted ($c(z, t = 0) = 0$), and suddenly brought into
contact with oxygen rich solvent at both the left and right side
$c_L = c_R = c_{\text{ref}}$ for $t \geq 0$. The concentration profile is divided
by the equilibrium profile, such that the black horizontal line
1 denotes equilibrium (to which the system converges for
larger times).

and described by averages for the permeant concentration
$c_M(t)$ and Boltzmann factor $e^{-\beta \Delta F_M}$.
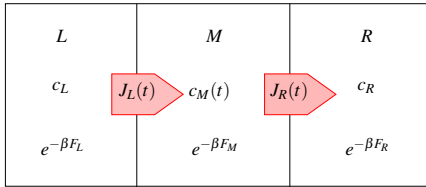
2. The phospholipid head-group regions ($z \in \{0, L\} \cup \{L -
h, h\}$) absorb no permeants, such that there exist steady-
state fluxes $J_L(t)$ and $J_R(t)$ through the left and right
head-group regions, respectively.

## S1.2 Three-state model

Using the above assumptions, the permeant kinetics can be described by the three states $L$ (left solvent region), $R$ (right solvent region), and $M$ (membrane), and the fluxes $J_L(t)$ and $J_R(t)$ at the solvent-membrane interfaces, as shown in Supp. Fig. 3. Each of these three states are defined by their (average) permeant concentration and (average) Boltzmann factor. The fluxes at the interfaces can be derived explicitly from the Smoluchowski equation, by imposing steady-state fluxes over the head-group regions and solving for $J_L(t)$ and $J_R(t)$

$$
\begin{aligned}
J_L(t) &= \frac{e^{\beta \Delta F(L)} c(L,t) - e^{\beta \Delta F_L} c_L(t)}{\int_0^L \frac{dz}{D(z) e^{-\beta \Delta F(z)}}} \\
&= \frac{e^{\beta \Delta F_M} c_M(t) - e^{\beta \Delta F_L} c_L(t)}{\int_0^L \frac{dz}{D(z) e^{-\beta \Delta F(z)}}} \\
&= \frac{\frac{S(t)}{C} - e^{\beta \Delta F_L} c_L(t)}{R_L}, \\
J_R(t) &= \frac{e^{\beta \Delta F_R} c_R(t) - e^{\beta \Delta F(L-h)} c(L-h,t)}{\int_{L-h}^h \frac{dz}{D(z) e^{-\beta \Delta F(z)}}} \\
&= \frac{e^{\beta \Delta F_R} c_R(t) - e^{\beta \Delta F_M} c_M(t)}{\int_{L-h}^h \frac{dz}{D(z) e^{-\beta \Delta F(z)}}} \\
&= \frac{e^{\beta \Delta F_R} c_R(t) - \frac{S(t)}{C}}{R_R} \\
&= \frac{e^{\beta \Delta F_R} c_R(t) - \frac{S(t)}{C}}{R_R}
\end{aligned} \tag{S6}
$$

In the second equalities, we have used the relation of Eq. S5 which relates $c_M(t)$ to $S(t)$ and $e^{-\beta \Delta F_M}$ to $C$ (or $Z$). We have also introduced the left and right head-group resistances $R_L$ and $R_R$, respectively. In the expressions for the fluxes $J_L(t)$ and $J_R(t)$, it is seen that the inversely Boltzmann weighted concentrations $c(z)/e^{-\beta F(z)}$ drive permeant displacement rather than just the concentration $c(z)$.



**Supplementary Figure 3.** Using the quasi-equilibrium approximation in the membrane interior and the steady-state approximation through the head-group regions, the charging process can be described using three states: (M) membrane, (L) left solvent region, and (R) right solvent region. The head-group regions are represented by the red arrows, through which the permeant fluxes $J_L(t)$ and $J_R(t)$ flow.

The differential equation is now found by integrating the

mass-balance law $\frac{\partial c(z,t)}{\partial t} = -\frac{\partial j(z,t)}{\partial z}$ over the membrane interior $z \in \{L, h - L\}$

$$
\frac{\partial}{\partial t} S(t) = -J_R(t) + J_L(t). \tag{S7}
$$

Plugging in the expressions of Eq. S6 results in the first order differential equation for $S(t)$

$$
\frac{\partial}{\partial t} S(t) = \frac{1}{C}\left(\frac{1}{R_L} + \frac{1}{R_R}\right) S(t) - \left(\frac{e^{\beta \Delta F_L} c_L(t)}{R_L} + \frac{e^{\beta \Delta F_R} c_R(t)}{R_R}\right), \tag{S8}
$$

where the time constant $\tau$ is given by

$$
\tau = \left[\frac{1}{C}\left(\frac{1}{R_L} + \frac{1}{R_R}\right)\right]^{-1}, \tag{S9}
$$

and the steady-state storage $S^{ss}$ is given by

$$
S^{ss} = \left(\frac{R_R}{R_R + R_L} e^{\beta \Delta F_L} c_L + \frac{R_L}{R_R + R_L} e^{\beta \Delta F_R} c_R\right) C, \tag{S10}
$$

where the latter equation is found by setting $\partial S/\partial t = 0$ and keeping $c_L$ and $c_R$ constant.

## S1.3 Discussion

First, if the membrane is symmetric, then we have $R_L = R_R = R_M/2$, and $\tau = (R_M C)/4$, which is, apart from smaller integration domains in the definitions of $R_L$ and $R_R$, equal to the time constant of the main manuscript. For a symmetric membrane ($R_L = R_R$) and equivalent solvents ($F_L = F_R$), the equilibrium ($c_L = c_R$) storage $S^{eq}$ (Eq. 4 in the main text) can be retrieved as follows. By setting $F_{ref} = F_L(= F_R)$ and denoting $c_{ref} = c_L(= c_R)$, we have that $\Delta F_L = \Delta F_R = 0$ and Eq. S10 becomes

$$
\begin{aligned}
S^{eq} &= c_{ref} C e^{\beta \times 0}, \\
&= c_{ref} \int_0^h e^{-\beta(F(z) - F_{ref})} dz. 
\end{aligned} \tag{S11}
$$

Second, we compare the time constant $\tau$ as predicted by our model with the time constant $\tau_R$ given by the discretized rate matrix $R$ for 81 systems, built by the 9 toy $F$ profiles and 9 toy $D$ profiles shown in Supp. Fig. 4. Excluding the 9 systems containing the flat free energy profile $F_8$ (for which our model should not hold, see bottom row of Supp. Fig. 4), the average ratio of $r = \tau/\tau_R$ is 0.96, with a standard deviation of 0.05.

Third, while the model accurately captures the time-constant, it does not yet perfectly capture the steady-state flux through the entire membrane. This is because the resistances $R_L$ and $R_R$ were only integrated over the phospholipid head-group regions, rather than over an entire leaflet. As the local resistance in the hydrophobic core is small compared to the local resistance of the head-group regions, one could extend the integration domain to include the membrane interior:
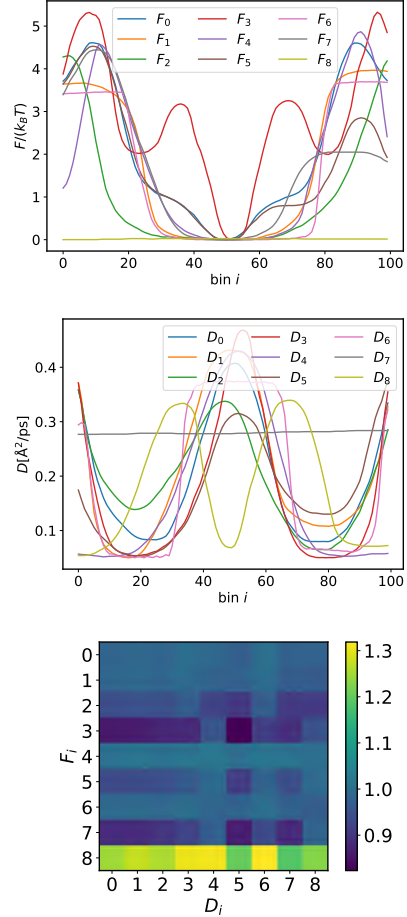
$$
\begin{aligned}
R_L &= \int_0^{h/2} \frac{dz}{D(z) e^{-\beta \Delta F(z)}}, \\
R_R &= \int_{h/2}^h \frac{dz}{D(z) e^{-\beta \Delta F(z)}},
\end{aligned} \tag{S12}
$$

where, for example, the integration domain for $R_L$ now runs from 0 to $h/2$, instead of 0 to $L$. In words: one could say that the permeant molecule has only entered state $M$ once it has reached the bilayer midplane. Doing so, we have that the total membrane resistance is given by $R_L + R_R$, and we have recovered the model that is described exactly by the RC network in the main text. Denoting the time-constants of this model by $\tau_{RC}$, we have that the average ratio $\tau_{RC}/\tau_R$ for the toy systems (excluding flat $F_8$ systems) is 1.04, with a standard deviation of 0.04.

Fourth, note that it are the leaflet resistances $R_L$ and $R_R$ that enter in the expression for $\tau$. If one were to set $R_L = R_R = R_M/2$, i.e. taking half of the total membrane resistance left and right, then the resulting $\tau_{\text{symmetric}} = (R_M C)/4$ has an average ratio $\tau_{\text{symmetric}}/\tau$ of 1.23 with standard deviation of 0.57 (excluding the systems with the flat $F_8$). The largest mismatch is for systems containing profiles $F_5$ or $F_7$ where the difference in $F_L$ and $F_R$ is no longer negligible.

Fifth, note that it are the leaflet resistances $R_L$ and $R_R$ that enter in the expression for the steady-state storage in Eq. S10.



**Supplementary Figure 4.** Nine free energy toy profiles (top) and nine diffusivity toy profiles (middle), where both consist of 100 bins. Each combination of an $F$ and $D$ profile represents a bilayer toy system. The ratio $r = \tau/\tau_R$ is shown on the bottom panel for each toy system. Note that the bottom row corresponds to a flat free energy surface, and therefore does not (and should not) satisfy the assumptions of our model.

# Bibliography

[1] J. Riskin, 'The defecating duck, or, the ambiguous origins of artificial life', *Critical inquiry*, vol. 29, no. 4, pp. 599–633, 2003.

[2] M. E. Moran, 'Jacques de vaucanson: The father of simulation', *Journal of endourology*, vol. 21, no. 7, pp. 679–683, 2007.

[3] H. M. Berman *et al.*, 'The protein data bank', *Nucleic acids research*, vol. 28, no. 1, pp. 235–242, 2000.

[4] L. M. Sampaleanu, F. Vallée, C. Slingsby and P. L. Howell, 'Structural studies of duck $\delta 1$ and $\delta 2$ crystallin suggest conformational changes occur during catalysis', *Biochemistry*, vol. 40, no. 9, pp. 2732–2742, 2001.

[5] R. Descartes, 'Meditations on first philosophy'. Broadview Press, 2013.

[6] F. Capra and P. L. Luisi, 'The mechanistic view of life', in *The Systems View of Life: A Unifying Vision*. Cambridge University Press, 2014, pp. 35–44.

[7] D. Jalobeanu, 'Constructing Natural Historical Facts: Baconian Natural History in Newton's First Paper on Light and Colors', in *Newton and Empiricism*, Oxford University Press, 2014.

[8] J. Louth, 'From newton to newtonianism: Reductionism and the development of the social sciences.' *Emergence: Complexity & Organization*, vol. 13, no. 4, 2011.

[9] W. D. Ross, 'Aristotle Metaphysics. A Revised Text with Introduction and Commentary'. Oxford: Oxford University Press, 1925.

[10] F. C. Fang and A. Casadevall, *Reductionistic and holistic science*, 2011.

[11] R. Dawkins, 'The Blind Watchmaker' (Penguin Books (Publisher). Penguin Science). Penguin, 2006.

[12]  A. Singharoy *et al.*, 'Atoms to phenotypes: Molecular design principles of cellular energy metabolism', *Cell*, vol. 179, no. 5, pp. 1098–1111, 2019.

[13]  M. E. Tuckerman, 'Statistical mechanics: theory and molecular simulation'. Oxford university press, 2023.

[14]  D. Frenkel and B. Smit, 'Understanding molecular simulation: from algorithms to applications'. Elsevier, 2023.

[15]  N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller and E. Teller, 'Equation of state calculations by fast computing machines', *The journal of chemical physics*, vol. 21, no. 6, pp. 1087–1092, 1953.

[16]  B. J. Alder, T. E. Wainwright *et al.*, 'Phase transition for a hard sphere system', *The Journal of chemical physics*, vol. 27, no. 5, p. 1208, 1957.

[17]  J. A. McCammon, B. R. Gelin and M. Karplus, 'Dynamics of folded proteins', *Nature*, vol. 267, no. 5612, pp. 585–590, 1977.

[18]  D. Macuglia, B. Roux and G. Ciccotti, 'The emergence of protein dynamics simulations: How computational statistical mechanics met biochemistry', *The European Physical Journal H*, vol. 47, no. 1, p. 13, 2022.

[19]  R. R. Schaller, 'Moore's law: Past, present and future', *IEEE spectrum*, vol. 34, no. 6, pp. 52–59, 1997.

[20]  J. Shalf, 'The future of computing beyond moore's law', *Philosophical Transactions of the Royal Society A*, vol. 378, no. 2166, p. 20 190 061, 2020.

[21]  P. L. Freddolino, A. S. Arkhipov, S. B. Larson, A. McPherson and K. Schulten, 'Molecular dynamics simulations of the complete satellite tobacco mosaic virus', *Structure*, vol. 14, no. 3, pp. 437–449, 2006.

[22]  J. Jung *et al.*, 'Scaling molecular dynamics beyond 100,000 processor cores for large-scale biophysical simulations', *Journal of computational chemistry*, vol. 40, no. 21, pp. 1919–1930, 2019.

[23]  L. Casalino *et al.*, 'Ai-driven multiscale simulations illuminate mechanisms of sars-cov-2 spike dynamics', *The International Journal of High Performance Computing Applications*, vol. 35, no. 5, pp. 432–451, 2021.

[24]   R. Milo, 'What is the total number of protein molecules per cell volume? a call to rethink some published values', *Bioessays*, vol. 35, no. 12, pp. 1050–1055, 2013.

[25]   R. J. Ellis, 'Macromolecular crowding: Obvious but under-appreciated', *Trends in biochemical sciences*, vol. 26, no. 10, pp. 597–604, 2001.

[26]   J. Zhang *et al.*, 'Artificial intelligence enhanced molecular simulations', *Journal of Chemical Theory and Computation*, vol. 19, no. 14, pp. 4338–4350, 2023.

[27]   C. Gupta, D. Sarkar, D. P. Tieleman and A. Singharoy, 'The ugly, bad, and good stories of large-scale biomolecular simulations', *Current Opinion in Structural Biology*, vol. 73, p. 102 338, 2022.

[28]   S. Y. Joshi and S. A. Deshmukh, 'A review of advancements in coarse-grained molecular dynamics simulations', *Molecular Simulation*, vol. 47, no. 10-11, pp. 786–803, 2021.

[29]   S. J. Marrink, H. J. Risselada, S. Yefimov, D. P. Tieleman and A. H. De Vries, 'The martini force field: Coarse grained model for biomolecular simulations', *The journal of physical chemistry B*, vol. 111, no. 27, pp. 7812–7824, 2007.

[30]   J. A. Stevens *et al.*, 'Molecular dynamics simulation of an entire cell', *Frontiers in Chemistry*, vol. 11, p. 1 106 495, 2023.

[31]   Y. Duan and P. A. Kollman, 'Pathways to a protein folding intermediate observed in a 1-microsecond simulation in aqueous solution', *Science*, vol. 282, no. 5389, pp. 740–744, 1998.

[32]   D. E. Shaw *et al.*, 'Millisecond-scale molecular dynamics simulations on anton', in *Proceedings of the conference on high performance computing networking, storage and analysis*, 2009, pp. 1–11.

[33]   D. E. Shaw *et al.*, 'Anton, a special-purpose machine for molecular dynamics simulation', *Commun. ACM*, vol. 51, no. 7, pp. 91–97, 2008.

[34]   D. E. Shaw *et al.*, 'Anton 2: Raising the bar for performance and programmability in a special-purpose molecular dynamics supercomputer', in *SC'14: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, IEEE, 2014, pp. 41–53.

[35]  D. E. Shaw *et al.*, 'Anton 3: Twenty microseconds of molecular dynamics simulation before lunch', in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, ser. SC '21, St. Louis, Missouri: Association for Computing Machinery, 2021.

[36]  A. Lyczek *et al.*, 'Mutation in abl kinase with altered drug-binding kinetics indicates a novel mechanism of imatinib resistance', *Proceedings of the National Academy of Sciences*, vol. 118, no. 46, e2111451118, 2021.

[37]  M. J. Abraham *et al.*, 'Gromacs: High performance molecular simulations through multi-level parallelism from laptops to supercomputers', *SoftwareX*, vol. 1, pp. 19–25, 2015.

[38]  L. Greengard and V. Rokhlin, 'A fast algorithm for particle simulations', *Journal of computational physics*, vol. 73, no. 2, pp. 325–348, 1987.

[39]  T. Darden, D. York and L. Pedersen, 'Particle mesh ewald: An n log n method for ewald sums in large systems', *The Journal of chemical physics*, vol. 98, no. 12, pp. 10 089–10 092, 1993.

[40]  J. Kurzak and B. M. Pettitt, 'Fast multipole methods for particle dynamics', *Molecular simulation*, vol. 32, no. 10-11, pp. 775–790, 2006.

[41]  D. Shamshirgar, R. Yokota, A.-K. Tornberg and B. Hess, 'Regularizing the fast multipole method for use in molecular simulation', *The Journal of Chemical Physics*, vol. 151, no. 23, 2019.

[42]  R. W. Hockney and J. W. Eastwood, 'Computer simulation using particles'. Boca Raton: CRC Press, 1988.

[43]  D. Brown, B. Maigret *et al.*, 'A domain decomposition parallel processing algorithm for molecular dynamics simulations of systems of arbitrary connectivity', *Computer Physics Communications*, vol. 103, no. 2-3, pp. 170–186, 1997.

[44]  M. P. Allen and D. J. Tildesley, 'Computer simulation of liquids'. Oxford university press, 2017.

[45]  J.-L. Lions, 'Résolution d'edp par un schéma en temps «pararéel» a"parareal"in time discretization of pde's', *Academie des Sciences Paris Comptes Rendus Serie Sciences Mathematiques*, vol. 332, no. 7, pp. 661–668, 2001.

[46] M. J. Gander and G. Wanner, 'From euler, ritz, and galerkin to modern computing', *Siam Review*, vol. 54, no. 4, pp. 627–666, 2012.

[47] R. D. Falgout, S. Friedhoff, T. V. Kolev, S. P. MacLachlan and J. B. Schroder, 'Parallel time integration with multigrid', *SIAM Journal on Scientific Computing*, vol. 36, no. 6, pp. C635–C661, 2014.

[48] B. Leimkuhler and C. Matthews, 'Molecular dynamics', *Interdisciplinary applied mathematics*, vol. 39, p. 443, 2015.

[49] M. Sahil, S. Sarkar and J. Mondal, 'Long-time-step molecular dynamics can retard simulation of protein-ligand recognition process', *Biophysical Journal*, vol. 122, no. 5, pp. 802–816, 2023.

[50] J.-P. Ryckaert, G. Ciccotti and H. J. Berendsen, 'Numerical integration of the cartesian equations of motion of a system with constraints: Molecular dynamics of n-alkanes', *Journal of computational physics*, vol. 23, no. 3, pp. 327–341, 1977.

[51] S. Miyamoto and P. A. Kollman, 'Settle: An analytical version of the shake and rattle algorithm for rigid water models', *Journal of computational chemistry*, vol. 13, no. 8, pp. 952–962, 1992.

[52] B. Hess, H. Bekker, H. J. Berendsen and J. G. Fraaije, 'Lincs: A linear constraint solver for molecular simulations', *Journal of computational chemistry*, vol. 18, no. 12, pp. 1463–1472, 1997.

[53] C. W. Hopkins, S. Le Grand, R. C. Walker and A. E. Roitberg, 'Long-time-step molecular dynamics through hydrogen mass repartitioning', *Journal of chemical theory and computation*, vol. 11, no. 4, pp. 1864–1874, 2015.

[54] K. A. Feenstra, B. Hess and H. J. Berendsen, 'Improving efficiency of large time-scale molecular dynamics simulations of hydrogen-rich systems', *Journal of computational chemistry*, vol. 20, no. 8, pp. 786–798, 1999.

[55] K. Olesen, N. Awasthi, D. S. Bruhn, W. Pezeshkian and H. Khandelia, 'Faster simulations with a 5 fs time step for lipids in the charmm force field', *Journal of chemical theory and computation*, vol. 14, no. 6, pp. 3342–3350, 2018.

[56] J. C. Phillips *et al.*, 'Scalable molecular dynamics with namd', *Journal of computational chemistry*, vol. 26, no. 16, pp. 1781–1802, 2005.

[57] D. A. Case *et al.*, 'Amber 10', 2008.

[58] B. R. Brooks *et al.*, 'Charmm: The biomolecular simulation program', *Journal of computational chemistry*, vol. 30, no. 10, pp. 1545–1614, 2009.

[59] M. J. Harvey, G. Giupponi and G. D. Fabritiis, 'Acemd: Accelerating biomolecular dynamics in the microsecond time scale', *Journal of chemical theory and computation*, vol. 5, no. 6, pp. 1632–1639, 2009.

[60] P. Eastman *et al.*, 'Openmm 7: Rapid development of high performance algorithms for molecular dynamics', *PLoS computational biology*, vol. 13, no. 7, e1005659, 2017.

[61] A. D. MacKerell Jr, 'Empirical force fields for biological macromolecules: Overview and issues', *Journal of computational chemistry*, vol. 25, no. 13, pp. 1584–1604, 2004.

[62] Z. Jing *et al.*, 'Polarizable force fields for biomolecular simulations: Recent advances and applications', *Annual Review of biophysics*, vol. 48, pp. 371–394, 2019.

[63] D. Van der Spoel, 'Systematic design of biomolecular force fields', *Current opinion in structural biology*, vol. 67, pp. 18–24, 2021.

[64] O. T. Unke *et al.*, 'Machine learning force fields', *Chemical Reviews*, vol. 121, no. 16, pp. 10 142–10 186, 2021.

[65] J. W. Carter, M. A. Gonzalez, N. J. Brooks, J. M. Seddon and F. Bresme, 'Flip-flop asymmetry of cholesterol in model membranes induced by thermal gradients', *Soft Matter*, vol. 16, no. 25, pp. 5925–5932, 2020.

[66] A. A. Gurtovenko and I. Vattulainen, 'Molecular mechanism for lipid flip-flops', *The Journal of Physical Chemistry B*, vol. 111, no. 48, pp. 13 554–13 559, 2007.

[67] V. Sharma and E. Mamontov, 'Multiscale lipid membrane dynamics as revealed by neutron spectroscopy', *Progress in Lipid Research*, vol. 87, p. 101 179, 2022.

[68] S. Qian, V. K. Sharma and L. A. Clifton, 'Understanding the structure and dynamics of complex biomembrane interactions by neutron scattering techniques', *Langmuir*, vol. 36, no. 50, pp. 15 189–15 211, 2020.

[69]  A. A. Kawale and B. M. Burmann, 'Characterization of back-bone dynamics using solution nmr spectroscopy to discern the functional plasticity of structurally analogous proteins', *STAR protocols*, vol. 2, no. 4, p. 100 919, 2021.

[70]  H. Ode, M. Nakashima, S. Kitamura, W. Sugiura and H. Sato, 'Molecular dynamics simulation in virus research', *Frontiers in microbiology*, vol. 3, p. 31 245, 2012.

[71]  K. Henzler-Wildman and D. Kern, 'Dynamic personalities of proteins', *Nature*, vol. 450, no. 7172, pp. 964–972, 2007.

[72]  K. H. Nam, *Molecular dynamics—from small molecules to macromolecules*, 2021.

[73]  A. M. Gomes, P. J. Costa and M. Machuqueiro, 'Recent advances on molecular dynamics-based techniques to address drug membrane permeability with atomistic detail', *BBA advances*, p. 100 099, 2023.

[74]  Y. Wang *et al.*, 'An experimentally validated approach to calculate the blood-brain barrier permeability of small molecules', *Scientific reports*, vol. 9, no. 1, p. 6117, 2019.

[75]  R. M. Venable, A. Kramer and R. W. Pastor, 'Molecular dynamics simulations of membrane permeability', *Chemical reviews*, vol. 119, no. 9, pp. 5954–5997, 2019.

[76]  B. Peters, 'Reaction rate theory and rare events'. Amsterdam: Elsevier, 2017.

[77]  F. Jensen, 'Activation energies and the arrhenius equation', *Quality and Reliability Engineering International*, vol. 1, no. 1, pp. 13–17, 1985.

[78]  L. Piela, J. Kostrowicki and H. A. Scheraga, 'On the multiple-minima problem in the conformational analysis of molecules: Deformation of the potential energy hypersurface by the diffusion equation method', *The Journal of Physical Chemistry*, vol. 93, no. 8, pp. 3339–3346, 1989.

[79]  H. Grubmüller, 'Predicting slow structural transitions in macromolecular systems: Conformational flooding', *Physical Review E*, vol. 52, no. 3, p. 2893, 1995.

[80]  T. Huber, A. E. Torda and W. F. Van Gunsteren, 'Local elevation: A method for improving the searching properties of molecular dynamics simulation', *Journal of computer-aided molecular design*, vol. 8, pp. 695–708, 1994.

[81]  A. Laio and M. Parrinello, 'Escaping free-energy minima', *Proceedings of the national academy of sciences*, vol. 99, no. 20, pp. 12 562–12 566, 2002.

[82]  A. Barducci, G. Bussi and M. Parrinello, 'Well-tempered metadynamics: A smoothly converging and tunable free-energy method', *Physical review letters*, vol. 100, no. 2, p. 020 603, 2008.

[83]  E. Darve and A. Pohorille, 'Calculating free energies using average force', *The Journal of chemical physics*, vol. 115, no. 20, pp. 9169–9183, 2001.

[84]  J. Comer, J. C. Gumbart, J. Hénin, T. Lelièvre, A. Pohorille and C. Chipot, 'The adaptive biasing force method: Everything you always wanted to know but were afraid to ask', *The Journal of Physical Chemistry B*, vol. 119, no. 3, pp. 1129–1151, 2015.

[85]  O. Valsson and M. Parrinello, 'Variational approach to enhanced sampling and free energy calculations', *Physical review letters*, vol. 113, no. 9, p. 090 601, 2014.

[86]  S. Kumar, J. M. Rosenberg, D. Bouzida, R. H. Swendsen and P. A. Kollman, 'The weighted histogram analysis method for free-energy calculations on biomolecules. i. the method', *Journal of computational chemistry*, vol. 13, no. 8, pp. 1011–1021, 1992.

[87]  M. R. Shirts and J. D. Chodera, 'Statistically optimal analysis of samples from multiple equilibrium states', *The Journal of chemical physics*, vol. 129, no. 12, 2008.

[88]  D. Hamelberg, J. Mongan and J. A. McCammon, 'Accelerated molecular dynamics: A promising and efficient simulation method for biomolecules', *The Journal of chemical physics*, vol. 120, no. 24, pp. 11 919–11 929, 2004.

[89]  D. J. Earl and M. W. Deem, 'Parallel tempering: Theory, applications, and new perspectives', *Physical Chemistry Chemical Physics*, vol. 7, no. 23, pp. 3910–3916, 2005.

[90]  R. H. Swendsen and J.-S. Wang, 'Replica monte carlo simulation of spin-glasses', *Physical review letters*, vol. 57, no. 21, p. 2607, 1986.

[91]  Y. Sugita and Y. Okamoto, 'Replica-exchange molecular dynamics method for protein folding', *Chemical physics letters*, vol. 314, no. 1-2, pp. 141–151, 1999.

[92]    E. Marinari and G. Parisi, 'Simulated tempering: A new monte carlo scheme', *Europhysics letters*, vol. 19, no. 6, p. 451, 1992.

[93]    F. Wang and D. P. Landau, 'Efficient, multiple-range random walk algorithm to calculate the density of states', *Physical review letters*, vol. 86, no. 10, p. 2050, 2001.

[94]    Y. Q. Gao, 'An integrate-over-temperature approach for enhanced sampling', *The Journal of chemical physics*, vol. 128, no. 6, 2008.

[95]    H. Eyring, 'The activated complex in chemical reactions', *The Journal of Chemical Physics*, vol. 3, no. 2, pp. 107–115, 1935.

[96]    E. Wigner, 'The transition state method', *Transactions of the Faraday Society*, vol. 34, pp. 29–41, 1938.

[97]    K. J. Laidler and M. C. King, 'The development of transition-state theory', *J. phys. Chem*, vol. 87, no. 15, pp. 2657–2664, 1983.

[98]    D. Chandler, 'Statistical mechanics of isomerization dynamics in liquids and the transition state approximation', *The Journal of Chemical Physics*, vol. 68, no. 6, pp. 2959–2970, 1978.

[99]    J. N. Onuchic, Z. Luthey-Schulten and P. G. Wolynes, 'Theory of protein folding: The energy landscape perspective', *Annual review of physical chemistry*, vol. 48, no. 1, pp. 545–600, 1997.

[100]   C. Dellago, P. G. Bolhuis, F. S. Csajka and D. Chandler, 'Transition path sampling and the calculation of rate constants', *The Journal of chemical physics*, vol. 108, no. 5, pp. 1964–1977, 1998.

[101]   C. Dellago, P. G. Bolhuis and D. Chandler, 'Efficient transition path sampling: Application to lennard-jones cluster rearrangements', *The Journal of chemical physics*, vol. 108, no. 22, pp. 9236–9245, 1998.

[102]   C. Dellago, P. G. Bolhuis and D. Chandler, 'On the calculation of reaction rate constants in the transition path ensemble', *The Journal of chemical physics*, vol. 110, no. 14, pp. 6617–6625, 1999.

[103]   P. G. Bolhuis, D. Chandler, C. Dellago and P. L. Geissler, 'Transition path sampling: Throwing ropes over rough mountain passes, in the dark', *Annual review of physical chemistry*, vol. 53, no. 1, pp. 291–318, 2002.

[104] C. Dellago, P. G. Bolhuis and P. L. Geissler, 'Transition path sampling', *Advances in chemical physics*, vol. 123, pp. 1–78, 2002.

[105] P. Bolhuis and C. Dellago, 'Practical and conceptual path sampling issues', *The European Physical Journal Special Topics*, vol. 224, no. 12, pp. 2409–2427, 2015.

[106] P. G. Bolhuis and D. W. Swenson, 'Transition path sampling as markov chain monte carlo of trajectories: Recent algorithms, software, applications, and future outlook', *Advanced Theory and Simulations*, vol. 4, no. 4, p. 2 000 237, 2021.

[107] J. Juraszek and P. G. Bolhuis, 'Sampling the multiple folding mechanisms of trp-cage in explicit solvent', *Proceedings of the National Academy of Sciences*, vol. 103, no. 43, pp. 15 859–15 864, 2006.

[108] J. Juraszek, J. Vreede and P. G. Bolhuis, 'Transition path sampling of protein conformational changes', *Chemical Physics*, vol. 396, pp. 30–44, 2012.

[109] R. B. Best and G. Hummer, 'Reaction coordinates and rates from transition paths', *Proceedings of the National Academy of Sciences*, vol. 102, no. 19, pp. 6732–6737, 2005.

[110] M. F. Hagan, A. R. Dinner, D. Chandler and A. K. Chakraborty, 'Atomistic understanding of kinetic pathways for single base-pair binding and unbinding in dna', *Proceedings of the National Academy of Sciences*, vol. 100, no. 24, pp. 13 922–13 927, 2003.

[111] J. Hu, A. Ma and A. R. Dinner, 'A two-step nucleotide-flipping mechanism enables kinetic discrimination of dna lesions by agt', *Proceedings of the National Academy of Sciences*, vol. 105, no. 12, pp. 4615–4620, 2008.

[112] T. S. Van Erp, D. Moroni and P. G. Bolhuis, 'A novel path sampling method for the calculation of rate constants', *The Journal of chemical physics*, vol. 118, no. 17, pp. 7762–7774, 2003.

[113] T. S. Van Erp, 'Dynamical rare event simulation techniques for equilibrium and nonequilibrium systems', *Advances in Chemical Physics*, vol. 151, p. 27, 2012.

[114] T. S. van Erp, 'How far can we stretch the timescale with retis?', *Europhysics Letters*, vol. 143, no. 3, p. 30 001, 2023.

[115]  T. S. van Erp, 'Reaction rate calculation by parallel path swapping', *Physical review letters*, vol. 98, no. 26, p. 268 301, 2007.

[116]  P. G. Bolhuis, 'Rare events via multiple reaction channels sampled by path replica exchange', *The Journal of chemical physics*, vol. 129, no. 11, 2008.

[117]  R. Cabriolu, K. M. Skjelbred Refsnes, P. G. Bolhuis and T. S. van Erp, 'Foundations and latest advances in replica exchange transition interface sampling', *The Journal of Chemical Physics*, vol. 147, no. 15, 2017.

[118]  S. Roet, D. T. Zhang and T. S. van Erp, 'Exchanging replicas with unequal cost, infinitely and permanently', *The Journal of Physical Chemistry A*, vol. 126, no. 47, pp. 8878–8886, 2022.

[119]  D. T. Zhang, L. Baldauf, S. Roet, A. Lervik and T. S. van Erp, 'Highly parallelizable path sampling with minimal rejections using asynchronous replica exchange and infinite swaps', *Proceedings of the National Academy of Sciences*, vol. 121, no. 7, e2318731121, 2024.

[120]  J. Rogal and P. G. Bolhuis, 'Multiple state transition path sampling', *The Journal of chemical physics*, vol. 129, no. 22, 2008.

[121]  D. Moroni, P. G. Bolhuis and T. S. van Erp, 'Rate constants for diffusive processes by partial path sampling', *The Journal of chemical physics*, vol. 120, no. 9, pp. 4055–4065, 2004.

[122]  W. Vervust, D. T. Zhang, T. S. Van Erp and A. Ghysels, 'Path sampling with memory reduction and replica exchange to reach long permeation timescales', *Biophysical Journal*, vol. 122, no. 14, pp. 2960–2972, 2023.

[123]  W. Vervust, D. T. Zhang, A. Ghysels, S. Roet, T. S. van Erp and E. Riccardi, 'Pyretis 3: Conquering rare and slow events without boundaries', *Journal of Computational Chemistry*, 2024.

[124]  M. Grünwald, C. Dellago and P. L. Geissler, 'Precision shooting: Sampling long transition pathways', *The Journal of chemical physics*, vol. 129, no. 19, 2008.

[125]  T. R. Gingrich and P. L. Geissler, 'Preserving correlations between trajectories for efficient path sampling', *The Journal of chemical physics*, vol. 142, no. 23, 2015.

[126] H. Jung, K.-i. Okazaki and G. Hummer, 'Transition path sampling of rare events by shooting from the top', *The Journal of chemical physics*, vol. 147, no. 15, 2017.

[127] G. Menzl, A. Singraber and C. Dellago, 'S-shooting: A bennett–chandler-like method for the computation of rate constants from committor trajectories', *Faraday Discussions*, vol. 195, pp. 345–364, 2016.

[128] E. E. Borrero and C. Dellago, 'Avoiding traps in trajectory space: Metadynamics enhanced transition path sampling', *The European Physical Journal Special Topics*, vol. 225, pp. 1609–1620, 2016.

[129] E. Riccardi, O. Dahlen and T. S. van Erp, 'Fast decorrelating monte carlo moves for efficient path sampling', *The Journal of Physical Chemistry Letters*, vol. 8, no. 18, pp. 4456–4460, 2017.

[130] D. T. Zhang, E. Riccardi and T. S. van Erp, 'Enhanced path sampling using subtrajectory monte carlo moves', *The Journal of Chemical Physics*, vol. 158, no. 2, 2023.

[131] R. J. Allen, P. B. Warren and P. R. Ten Wolde, 'Sampling rare switching events in biochemical networks', *Physical review letters*, vol. 94, no. 1, p. 018104, 2005.

[132] P. Melnik-Melnikov and E. Dekhtyaruk, 'Rare events probabilities estimation by "russian roulette and splitting" simulation technique', *Probabilistic engineering mechanics*, vol. 15, no. 2, pp. 125–129, 2000.

[133] T. S. van Erp, 'Efficiency analysis of reaction rate calculation methods using analytical models i: The two-dimensional sharp barrier', *The Journal of chemical physics*, vol. 125, no. 17, 2006.

[134] M. Villén-Altamirano and J. Villen-Altamirano, 'Analysis of restart simulation: Theoretical basis and sensitivity study', *European Transactions on Telecommunications*, vol. 13, no. 4, pp. 373–385, 2002.

[135] 'Weighted-ensemble brownian dynamics simulations for protein association reactions', *Biophysical journal*, vol. 70, no. 1, pp. 97–110, 1996.

[136] T. Sztain *et al.*, 'A glycan gate controls opening of the sars-cov-2 spike protein', *Nature chemistry*, vol. 13, no. 10, pp. 963–968, 2021.

[137]  F. Cérou, A. Guyader, T. Lelievre and D. Pommier, 'A multiple replica approach to simulate reactive trajectories', *The Journal of chemical physics*, vol. 134, no. 5, 2011.

[138]  A. K. Faradjian and R. Elber, 'Computing time scales from reaction coordinates by milestoning', *The Journal of chemical physics*, vol. 120, no. 23, pp. 10 880–10 889, 2004.

[139]  E. Vanden-Eijnden and M. Venturoli, 'Markovian milestoning with voronoi tessellations', *The Journal of chemical physics*, vol. 130, no. 19, 2009.

[140]  P. Májek and R. Elber, 'Milestoning without a reaction coordinate', *Journal of chemical theory and computation*, vol. 6, no. 6, pp. 1805–1817, 2010.

[141]  S. Kirmizialtin and R. Elber, 'Revisiting and computing reaction coordinates with directional milestoning', *The journal of physical chemistry A*, vol. 115, no. 23, pp. 6137–6148, 2011.

[142]  J. M. Bello-Rivas and R. Elber, 'Exact milestoning', *The Journal of Chemical Physics*, vol. 142, no. 9, 2015.

[143]  R. Elber, 'Milestoning: An efficient approach for atomically detailed simulations of kinetics in biophysics', *Annual review of biophysics*, vol. 49, pp. 69–85, 2020.

[144]  A. M. Berezhkovskii and A. Szabo, 'Committors, first-passage times, fluxes, markov states, milestones, and all that', *The Journal of chemical physics*, vol. 150, no. 5, 2019.

[145]  N. S. Hinrichs and V. S. Pande, 'Calculation of the distribution of eigenvalues and eigenvectors in markovian state models for molecular dynamics', *The Journal of chemical physics*, vol. 126, no. 24, 2007.

[146]  N.-V. Buchete and G. Hummer, 'Coarse master equations for peptide folding dynamics', *The Journal of Physical Chemistry B*, vol. 112, no. 19, pp. 6057–6069, 2008.

[147]  M. Sarich, F. Noé and C. Schütte, 'On the approximation quality of markov state models', *Multiscale Modeling & Simulation*, vol. 8, no. 4, pp. 1154–1177, 2010.

[148]  G. R. Bowman, V. S. Pande and F. Noé, 'An introduction to Markov state models and their application to long timescale molecular simulation'. Springer Science & Business Media, 2013, vol. 797.

[149] B. E. Husic and V. S. Pande, 'Markov state models: From an art to a science', *Journal of the American Chemical Society*, vol. 140, no. 7, pp. 2386–2396, 2018.

[150] P. Deuflhard, W. Huisinga, A. Fischer and C. Schütte, 'Identification of almost invariant aggregates in reversible nearly uncoupled markov chains', *Linear Algebra and its Applications*, vol. 315, no. 1-3, pp. 39–59, 2000.

[151] P. Deuflhard and M. Weber, 'Robust perron cluster analysis in conformation dynamics', *Linear algebra and its applications*, vol. 398, pp. 161–184, 2005.

[152] G. R. Bowman, 'Improved coarse-graining of markov state models via explicit consideration of statistical uncertainty', *The Journal of Chemical Physics*, vol. 137, no. 13, 2012.

[153] B. E. Husic, K. A. McKiernan, H. K. Wayment-Steele, M. M. Sultan and V. S. Pande, 'A minimum variance clustering approach produces robust and interpretable coarse-grained models', *Journal of chemical theory and computation*, vol. 14, no. 2, pp. 1071–1082, 2018.

[154] F. Noé and F. Nuske, 'A variational approach to modeling slow processes in stochastic dynamical systems', *Multiscale Modeling & Simulation*, vol. 11, no. 2, pp. 635–655, 2013.

[155] Y. Naritomi and S. Fuchigami, 'Slow dynamics in protein fluctuations revealed by time-structure based independent component analysis: The case of domain motions', *The Journal of chemical physics*, vol. 134, no. 6, 2011.

[156] C. Schütte, S. Klus and C. Hartmann, 'Overcoming the timescale barrier in molecular dynamics: Transfer operators, variational principles and machine learning', *Acta Numerica*, vol. 32, pp. 517–673, 2023.

[157] A. Mardt, L. Pasquali, H. Wu and F. Noé, 'Vampnets for deep learning of molecular kinetics', *Nature communications*, vol. 9, no. 1, p. 5, 2018.

[158] C. A. Lipinski, F. Lombardo, B. W. Dominy and P. J. Feeney, 'Experimental and computational approaches to estimate solubility and permeability in drug discovery and development settings', *Advanced drug delivery reviews*, vol. 64, pp. 4–17, 2012.

[159] W. L. Jorgensen, 'The many roles of computation in drug discovery', *Science*, vol. 303, no. 5665, pp. 1813–1818, 2004.

[160] R. Claveria-Gimeno, S. Vega, O. Abian and A. Velazquez-Campoy, 'A look at ligand binding thermodynamics in drug discovery', *Expert Opinion on Drug Discovery*, vol. 12, no. 4, pp. 363–377, 2017.

[161] D. L. Mobley and M. K. Gilson, 'Predicting binding free energies: Frontiers and benchmarks', *Annual review of biophysics*, vol. 46, pp. 531–558, 2017.

[162] R. A. Copeland, D. L. Pompliano and T. D. Meek, 'Drug–target residence time and its implications for lead optimization', *Nature reviews Drug discovery*, vol. 5, no. 9, pp. 730–739, 2006.

[163] P. J. Tummino and R. A. Copeland, 'Residence time of receptor- ligand complexes and its effect on biological function', *Biochemistry*, vol. 47, no. 20, pp. 5481–5492, 2008.

[164] R. A. Copeland, 'The dynamics of drug-target interactions: Drug-target residence time and its impact on efficacy and safety', *Expert opinion on drug discovery*, vol. 5, no. 4, pp. 305–310, 2010.

[165] H. Lu and P. J. Tonge, 'Drug–target residence time: Critical information for lead optimization', *Current opinion in chemical biology*, vol. 14, no. 4, pp. 467–474, 2010.

[166] R. A. Copeland, 'The drug–target residence time model: A 10-year retrospective', *Nature Reviews Drug Discovery*, vol. 15, no. 2, pp. 87–95, 2016.

[167] B. K. Shoichet, 'Virtual screening of chemical libraries', *Nature*, vol. 432, no. 7019, pp. 862–865, 2004.

[168] N. Okimoto *et al.*, 'High-performance drug discovery: Computational screening by combining docking and molecular dynamics simulations', *PLoS computational biology*, vol. 5, no. 10, e1000528, 2009.

[169] W. L. Jorgensen, 'Efficient drug lead discovery and optimization', *Accounts of chemical research*, vol. 42, no. 6, pp. 724–733, 2009.

[170] X. Lin, X. Li and X. Lin, 'A review on applications of computational methods in drug screening and design', *Molecules*, vol. 25, no. 6, p. 1375, 2020.

[171] L. Wang, H. L. McLeod and R. M. Weinshilboum, 'Genomics and drug response', *New England Journal of Medicine*, vol. 364, no. 12, pp. 1144–1153, 2011.

[172]  M. V. Relling and W. E. Evans, 'Pharmacogenomics in the clinic', *Nature*, vol. 526, no. 7573, pp. 343–350, 2015.

[173]  P. Sneha and C. George Priya Doss, 'Chapter seven - molecular dynamics: New frontier in personalized medicine', in *Personalized Medicine*, ser. Advances in Protein Chemistry and Structural Biology, R. Donev, Ed., vol. 102, Academic Press, 2016, pp. 181–224.

[174]  F. Hyder, D. L. Rothman and M. R. Bennett, 'Cortical energy demands of signaling and nonsignaling components in brain are conserved across mammalian species and activity levels', *Proceedings of the National Academy of Sciences*, vol. 110, no. 9, pp. 3549–3554, 2013.

[175]  M. E. Watts, R. Pocock and C. Claudianos, 'Brain energy and oxygen metabolism: Emerging role in normal function and disease', *Frontiers in molecular neuroscience*, vol. 11, p. 216, 2018.

[176]  K.-A. Nave and H. B. Werner, 'Myelination of the nervous system: Mechanisms and functions', *Annual review of cell and developmental biology*, vol. 30, pp. 503–533, 2014.

[177]  C. Dellago, P. G. Bolhuis and P. L. Geissler, 'Transition path sampling methods', *Computer Simulations in Condensed Matter Systems: From Materials to Chemical Biology Volume 1*, pp. 349–391, 2006.

[178]  K. Lindorff-Larsen, P. Maragakis, S. Piana, M. P. Eastwood, R. O. Dror and D. E. Shaw, 'Systematic validation of protein force fields against experimental data', *PloS one*, vol. 7, no. 2, e32131, 2012.

[179]  P. Dauber-Osguthorpe and A. T. Hagler, 'Biomolecular force fields: Where have we been, where are we now, where do we need to go and how do we get there?', *Journal of computer-aided molecular design*, vol. 33, no. 2, pp. 133–203, 2019.

[180]  B. Peters and B. L. Trout, 'Obtaining reaction coordinates by likelihood maximization', *The Journal of chemical physics*, vol. 125, no. 5, 2006.

[181]  J. Y. Wang, 'The capable abl: What is its biological function?', *Molecular and cellular biology*, 2014.

[182]  P. S. Georgoulia, G. Todde, S. Bjelic and R. Friedman, 'The catalytic activity of abl1 single and compound mutations: Implications for the mechanism of drug resistance mutations in chronic myeloid leukaemia', *Biochimica et Biophysica Acta (BBA)-General Subjects*, vol. 1863, no. 4, pp. 732–741, 2019.

[183]  S. Panjarian, R. E. Iacob, S. Chen, J. R. Engen and T. E. Smithgall, 'Structure and dynamic regulation of abl kinases', *Journal of Biological Chemistry*, vol. 288, no. 8, pp. 5443–5450, 2013.

[184]  T. O'hare, M. S. Zabriskie, A. M. Eiring and M. W. Deininger, 'Pushing the limits of targeted therapy in chronic myeloid leukaemia', *Nature Reviews Cancer*, vol. 12, no. 8, pp. 513–526, 2012.

[185]  F. Stegmeier, M. Warmuth, W. Sellers and M. Dorsch, 'Targeted cancer therapies in the twenty-first century: Lessons from imatinib', *Clinical Pharmacology & Therapeutics*, vol. 87, no. 5, pp. 543–552, 2010.

[186]  E. P. Reddy and A. K. Aggarwal, 'The ins and outs of bcr-abl inhibition', *Genes & cancer*, vol. 3, no. 5-6, pp. 447–454, 2012.

[187]  G. L. Simpson *et al.*, 'Identification and optimization of novel small c-abl kinase activators using fragment and hts methodologies', *Journal of medicinal chemistry*, vol. 62, no. 4, pp. 2154–2171, 2019.

[188]  T. Xie, T. Saleh, P. Rossi and C. G. Kalodimos, 'Conformational states dynamically populated by a kinase determine its function', *Science*, vol. 370, no. 6513, eabc2754, 2020.

[189]  W. L. DeLano *et al.*, 'Pymol: An open-source molecular graphics tool', *CCP4 Newsl. Protein Crystallogr*, vol. 40, no. 1, pp. 82–92, 2002.

[190]  S. Jo, T. Kim, V. G. Iyer and W. Im, 'Charmm-gui: A web-based graphical user interface for charmm', *Journal of computational chemistry*, vol. 29, no. 11, pp. 1859–1865, 2008.

[191]  W. W. Chan *et al.*, 'Conformational control inhibition of the bcr-abl1 tyrosine kinase, including the gatekeeper t315i mutant, by the switch-control inhibitor dcc-2036', *Cancer cell*, vol. 19, no. 4, pp. 556–568, 2011.

[192]  W. Humphrey, A. Dalke and K. Schulten, 'Vmd: Visual molecular dynamics', *Journal of molecular graphics*, vol. 14, no. 1, pp. 33–38, 1996.

[193]  J. Huang *et al.*, 'Charmm36m: An improved force field for folded and intrinsically disordered proteins', *Nature methods*, vol. 14, no. 1, pp. 71–73, 2017.

[194] W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey and M. L. Klein, 'Comparison of simple potential functions for simulating liquid water', *The Journal of chemical physics*, vol. 79, no. 2, pp. 926–935, 1983.

[195] Y.-L. Lin, Y. Meng, L. Huang and B. Roux, 'Computational study of gleevec and g6g reveals molecular determinants of kinase inhibitor selectivity', *Journal of the American Chemical Society*, vol. 136, no. 42, pp. 14 753–14 762, 2014.

[196] L. Huang and B. Roux, 'Automated force field parameterization for nonpolarizable and polarizable atomic models based on ab initio target data', *Journal of chemical theory and computation*, vol. 9, no. 8, pp. 3543–3556, 2013.

[197] E. Boulanger, L. Huang, C. Rupakheti, A. D. MacKerell Jr and B. Roux, 'Optimized lennard-jones parameters for druglike small molecules', *Journal of chemical theory and computation*, vol. 14, no. 6, pp. 3121–3131, 2018.

[198] F. Paul, T. Thomas and B. Roux, 'Diversity of long-lived intermediates along the binding pathway of imatinib to abl kinase revealed by md simulations', *Journal of chemical theory and computation*, vol. 16, no. 12, pp. 7852–7865, 2020.

[199] A. Aleksandrov and T. Simonson, 'A molecular mechanics model for imatinib and imatinib: Kinase binding', *Journal of computational chemistry*, vol. 31, no. 7, pp. 1550–1560, 2010.

[200] G. Bussi, D. Donadio and M. Parrinello, 'Canonical sampling through velocity rescaling', *The Journal of chemical physics*, vol. 126, no. 1, 2007.

[201] M. Parrinello and A. Rahman, 'Polymorphic transitions in single crystals: A new molecular dynamics method', *Journal of Applied physics*, vol. 52, no. 12, pp. 7182–7190, 1981.

[202] G. A. Tribello, M. Bonomi, D. Branduardi, C. Camilloni and G. Bussi, 'Plumed 2: New feathers for an old bird', *Computer physics communications*, vol. 185, no. 2, pp. 604–613, 2014.

[203] J. S. Hub, B. L. De Groot and D. Van Der Spoel, 'Gwham a free weighted histogram analysis implementation including robust error and autocorrelation estimates', *Journal of chemical theory and computation*, vol. 6, no. 12, pp. 3713–3720, 2010.

[204] R. T. McGibbon *et al.*, 'Mdtraj: A modern open library for the analysis of molecular dynamics trajectories', *Biophysical journal*, vol. 109, no. 8, pp. 1528–1532, 2015.

[205] N. Michaud-Agrawal, E. J. Denning, T. B. Woolf and O. Beckstein, 'Mdanalysis: A toolkit for the analysis of molecular dynamics simulations', *Journal of computational chemistry*, vol. 32, no. 10, pp. 2319–2327, 2011.

[206] M. A. Seeliger, B. Nagar, F. Frank, X. Cao, M. N. Henderson and J. Kuriyan, 'C-src binds to the cancer drug imatinib with an inactive abl/c-kit conformation and a distributed thermodynamic penalty', *Structure*, vol. 15, no. 3, pp. 299–311, 2007.

[207] R. V. Agafonov, C. Wilson, R. Otten, V. Buosi and D. Kern, 'Energetic dissection of gleevec's selectivity toward human tyrosine kinases', *Nature structural & molecular biology*, vol. 21, no. 10, pp. 848–853, 2014.

[208] B. Narayan, N.-V. Buchete and R. Elber, 'Computer simulations of the dissociation mechanism of gleevec from abl kinase with milestoning', *The Journal of Physical Chemistry B*, vol. 125, no. 22, pp. 5706–5715, 2021.

[209] M. Shekhar, Z. Smith, M. A. Seeliger and P. Tiwary, 'Protein flexibility and dissociation pathway differentiation can explain onset of resistance mutations in kinases', *Angewandte Chemie International Edition*, vol. 61, no. 28, e202200983, 2022.

[210] P. Ayaz *et al.*, 'Structural mechanism of a drug-binding process involving a large conformational change of the protein target', *Nature Communications*, vol. 14, no. 1, p. 1885, 2023.

[211] L. Verlet, 'Computer experiments on classical fluids. I. Thermodynamical properties of Lennard-Jones molecules', *Phys. Rev.*, vol. 159, pp. 98–103, 1 1967.

[212] M. P. Murphy, 'How mitochondria produce reactive oxygen species', *Biochemical journal*, vol. 417, no. 1, pp. 1–13, 2009.

[213] M. Erecińska and I. A. Silver, 'Tissue oxygen tension and brain sensitivity to hypoxia', *Respiration Physiology*, vol. 128, no. 3, pp. 263–276, 2001.

[214] K. Krab, H. Kempe and M. Wikström, 'Explaining the enigmatic km for oxygen in cytochrome c oxidase: A kinetic model', *Biochimica et Biophysica Acta (BBA) - Bioenergetics*, vol. 1807, no. 3, pp. 348–358, 2011.

[215] E. Gnaiger, B. Lassnig, A. Kuznetsov, G. Rieger and R. Margreiter, 'Mitochondrial Oxygen Affinity, Respiratory Flux Control And Excess Capacity Of Cytochrome c Oxidase', *Journal of Experimental Biology*, vol. 201, no. 8, pp. 1129–1139, 1998.

[216] M. Simons and K.-A. Nave, 'Oligodendrocytes: Myelination and axonal support', *Cold Spring Harbor perspectives in biology*, vol. 8, no. 1, a020479, 2016.

[217] A. Kister and I. Kister, 'Overview of myelin, major myelin lipids, and myelin-associated proteins', *Frontiers in Chemistry*, vol. 10, p. 1 041 961, 2023.

[218] A. Missner and P. Pohl, '110 years of the meyer–overton rule: Predicting membrane permeability of gases and other small compounds', *ChemPhysChem*, vol. 10, no. 9-10, pp. 1405–1414, 2009.

[219] A. Ghysels, R. M. Venable, R. W. Pastor and G. Hummer, 'Position-dependent diffusion tensors in anisotropic media from simulation: Oxygen transport in and through membranes', *Journal of chemical theory and computation*, vol. 13, no. 6, pp. 2962–2976, 2017.

[220] E. Riccardi, A. Krämer, T. S. van Erp and A. Ghysels, 'Permeation rates of oxygen through a lipid bilayer using replica exchange transition interface sampling', *The Journal of Physical Chemistry B*, vol. 125, no. 1, pp. 193–201, 2020.

[221] O. De Vos, R. M. Venable, T. Van Hecke, G. Hummer, R. W. Pastor and A. Ghysels, 'Membrane permeability: Characteristic times and lengths for oxygen and a simulation-based test of the inhomogeneous solubility-diffusion model', *Journal of chemical theory and computation*, vol. 14, no. 7, pp. 3811–3824, 2018.

[222] S. C. Pias, 'How does oxygen diffuse from capillaries to tissue mitochondria? barriers and pathways', *The Journal of Physiology*, vol. 599, no. 6, pp. 1769–1782, 2021.

[223] M. Moller, H. Botti, C. Batthyany, H. Rubbo, R. Radi and A. Denicola, 'Direct measurement of nitric oxide and oxygen partitioning into liposomes and low density lipoprotein', *Journal of Biological Chemistry*, vol. 280, no. 10, pp. 8850–8854, 2005.

[224] M. N. Möller, Q. Li, M. Chinnaraj, H. C. Cheung, J. R. Lancaster Jr and A. Denicola, 'Solubility and diffusion of oxygen in phospholipid membranes', *Biochimica et biophysica acta (bba)-biomembranes*, vol. 1858, no. 11, pp. 2923–2930, 2016.

[225] M. N. Möller and A. Denicola, 'Diffusion of nitric oxide and oxygen in lipoproteins and membranes studied by pyrene fluorescence quenching', *Free Radical Biology and Medicine*, vol. 128, pp. 137–143, 2018.

[226] M. S. Al-Abdul-Wahid, C.-H. Yu, I. Batruch, F. Evanics, R. Pomès and R. S. Prosser, 'A combined nmr and molecular dynamics study of the transmembrane solubility and diffusion rate profile of dioxygen in lipid bilayers', *Biochemistry*, vol. 45, no. 35, pp. 10 719–10 728, 2006.

[227] G. Hummer, 'Position-dependent diffusion coefficients and free energies from bayesian analysis of equilibrium and replica molecular dynamics simulations', *New Journal of Physics*, vol. 7, no. 1, p. 34, 2005.

[228] R. B. Buxton, K. Uludağ, D. J. Dubowitz and T. T. Liu, 'Modeling the hemodynamic response to brain activation', *NeuroImage*, vol. 23, S220–S233, 2004, Mathematics in Brain Imaging.

[229] R. Buxton, 'Interpreting oxygenation-based neuroimaging signals: The importance and the challenge of understanding brain oxygen metabolism', *Frontiers in neuroenergetics*, vol. 2, p. 1648, 2010.

[230] E. M. Hillman, 'Coupling mechanism and significance of the bold signal: A status report', *Annual review of neuroscience*, vol. 37, pp. 161–181, 2014.

[231] A. Devor, A. K. Dunn, M. L. Andermann, I. Ulbert, D. A. Boas and A. M. Dale, 'Coupling of total hemoglobin concentration, oxygenation, and neural activity in rat somatosensory cortex', *Neuron*, vol. 39, no. 2, pp. 353–359, 2003.

[232] A. Parpaleix, Y. G. Houssen and S. Charpak, 'Imaging local neuronal activity by monitoring po2 transients in capillaries', *Nature medicine*, vol. 19, no. 2, pp. 241–246, 2013.

[233] S. Sakadžić *et al.*, 'Two-photon microscopy measurement of cerebral metabolic rate of oxygen using periarteriolar oxygen concentration gradients', *Neurophotonics*, vol. 3, no. 4, pp. 045 005–045 005, 2016.

[234] J. Lecoq, P. Tiret, M. Najac, G. M. Shepherd, C. A. Greer and S. Charpak, 'Odor-evoked oxygen consumption by action potential and synaptic transmission in the olfactory bulb', *Journal of Neuroscience*, vol. 29, no. 5, pp. 1424–1433, 2009.

[235] H. S. Wei *et al.*, 'Erythrocytes are oxygen-sensing regulators of the cerebral microcirculation', *Neuron*, vol. 91, no. 4, pp. 851–862, 2016.

[236] C. Huchzermeyer, N. Berndt, H.-G. Holzhütter and O. Kann, 'Oxygen consumption rates during three different neuronal activity states in the hippocampal ca3 network', *Journal of Cerebral Blood Flow & Metabolism*, vol. 33, no. 2, pp. 263–271, 2013.

[237] J. Jumper *et al.*, 'Highly accurate protein structure prediction with alphafold', *Nature*, vol. 596, no. 7873, pp. 583–589, 2021.

[238] J. Abramson *et al.*, 'Accurate structure prediction of biomolecular interactions with alphafold 3', *Nature*, pp. 1–3, 2024.

[239] T. B. Brown *et al.*, *Language models are few-shot learners*, 2020.

[240] A. Gaulton *et al.*, 'The chembl database in 2017', *Nucleic acids research*, vol. 45, no. D1, pp. D945–D954, 2017.

[241] S. Kim *et al.*, 'Pubchem 2023 update', *Nucleic acids research*, vol. 51, no. D1, pp. D1373–D1380, 2023.

[242] M. K. Gilson, T. Liu, M. Baitaluk, G. Nicola, L. Hwang and J. Chong, 'Bindingdb in 2015: A public database for medicinal chemistry, computational chemistry and systems pharmacology', *Nucleic acids research*, vol. 44, no. D1, pp. D1045–D1053, 2016.

[243] C. Knox *et al.*, 'Drugbank 6.0: The drugbank knowledgebase for 2024', *Nucleic acids research*, vol. 52, no. D1, pp. D1265–D1275, 2024.

[244] H. M. Berman *et al.*, 'The protein data bank', *Acta Crystallographica Section D: Biological Crystallography*, vol. 58, no. 6, pp. 899–907, 2002.

[245] U. Consortium, 'Uniprot: A worldwide hub of protein knowledge', *Nucleic acids research*, vol. 47, no. D1, pp. D506–D515, 2019.

[246] D. A. Schuetz *et al.*, 'Kinetics for drug discovery: An industry-driven effort to target drug residence time', *Drug discovery today*, vol. 22, no. 6, pp. 896–911, 2017.

[247] H. Jung *et al.*, 'Machine-guided path sampling to discover mechanisms of molecular self-organization', *Nature Computational Science*, vol. 3, no. 4, pp. 334–345, 2023.