

Ab Initio Molecular Dynamics Simulations of Phosphate Hydrolysis Using Neural Network Potentials

Albert MAKHMUDOV

Supervisor: Prof. J. Harvey
KU Leuven

Thesis presented in
fulfillment of the requirements
for the degree of Master of Science
in Theoretical Chemistry and Computational Modelling

Academic year 2024-2025

© Copyright by KU Leuven

Without written permission of the promtors and the authors it is forbidden to reproduce or adapt in any form or by any means any part of this publication. Requests for obtaining the right to reproduce or utilize parts of this publication should be addressed to KU Leuven, Faculteit Wetenschappen, Celestijnenlaan 200H bus 2100, 3001 Leuven (Heverlee), telephone +32 16 32 14 01.

A written permission of the promotor is also required to use the methods, products, schematics and programs described in this work for industrial or commercial use, and for submitting this publication in scientific contests.

This thesis is an exam document that obtained no further correction of possible errors after the defense. Referring to this thesis in papers and analogous documents is only allowed after written consent of the supervisor(s), mentioned on the title page.

Foreword

Contribution statement

Summary

List of abbreviations

ADP adenosine diphosphate

ATP adenosine triphosphate

CN coordination number

CV collective variable

DNA deoxyribonucleic acid

GPW Gaussian plane wave method

GTH Goedecker-Teter-Hutter pseudopotentials

lg leaving group

MeDP methyl diphosphate trianion

MeHDP methyl diphosphate dianion

MeHMP methyl phosphate monoanion

MeMP methyl phosphate dianion

MeTP methyl triphosphate tetraniion

NNP neural network potential

NPT constant number of particles, pressure and temperature

Nuc nucleophile

NVT constant number of particles, volume and temperature

PBC periodic boundary conditions

PBE Perdew-Burke-Ernzerhof exchange-correlation functional

P_i inorganic phosphate

PME Particle Mesh Ewald

polyP polyphosphate

PT proton transfer

RMSE root mean square error

RNA ribonucleic acid

RRHO Rigid-Rotor Harmonic-Oscillator

SCF self-consistent field

TS transition state

TZV2P triple- ζ valence basis set with two polarisation functions

WTMetaD well-tempered metadynamics

MFEP minimum free energy path

Contents

1	Introduction	1
1.1	Role of phosphates in biological systems	1
1.2	Enzymes involved in phosphate hydrolysis	3
1.3	Reaction mechanism: phosphate monoesters	6
1.3.1	Phosphates	6
1.3.2	Diphosphates	8
1.3.3	Triphosphates	12
1.4	Research goals	12
2	Theory	14
2.1	A brief introduction to statistical mechanics	14
2.1.1	Partition functions	14
2.1.2	Macroscopic properties and thermodynamic functions	16
2.1.3	The canonical ensemble	17
2.1.4	Classical forcefields and molecular dynamics	18
2.1.5	Enhanced sampling techniques	18
2.2	Transition state theory	19
2.3	Density functional theory	19
2.3.1	The Kohn-Sham approach	19
2.3.2	Generalised gradient approximation and PBE functional	19
2.3.3	<i>Ab initio</i> molecular dynamics and GPW method	19
2.4	Extended tight binding	19
2.5	Neural network potentials	19
2.5.1	Message passing graph neural networks	19
2.5.2	Invariance and equivariance	19
2.5.3	Equivariant neural network potentials	19
3	Computational details	20
3.1	Training dataset generation	20
3.1.1	System preparation	20

3.1.2	Initial equilibration using classical force fields	21
3.1.3	Collective variables	22
3.1.4	GFN1-xTB based exploration of the configuration space	24
3.1.5	Data labeling	24
3.1.6	Iterative training of the neural network potential	25
3.2	Production runs at different temperatures	28
3.3	Validation of the transition states	29
3.4	Lifetime of the transition states	29
3.5	Data analysis and visualisation	29
4	Results and discussion	30
4.1	Final dataset composition	30
4.2	Accuracy of the neural network potential	31
4.3	Performance of the potential	33
4.4	Stability of the production runs	34
4.5	Convergence of the free energy profiles	34
4.6	Evolution of the collective variables over time	34
4.7	Reaction mechanism in case of the methyl diphosphate trianion	34
4.7.1	Minimum free energy path and free energy surface	34
4.7.2	Proton transfer mechanism	34
4.8	Reaction mechanism in case of the methyl diphosphate dianion	34
4.8.1	Minimum free energy path and free energy surface	34
4.8.2	Proton transfer mechanism	34
4.9	Arrhenius relationship	34
5	Conclusions	35
Bibliography		36
A Supplementary information		41

Chapter 1

Introduction

This chapter provides an overview of the role of phosphates in biological systems, the enzymes involved in phosphate hydrolysis, and a detailed explanation of the reaction mechanisms associated with these processes—topics that have puzzled researchers for a long time. The discussion begins with the fundamental importance of phosphates in life, particularly in energy transfer and storage. This is followed by a brief overview of the enzymes that catalyse phosphate hydrolysis and their implications for cellular function. Finally, the chapter explores the reaction mechanisms of phosphate hydrolysis, highlighting key studies and findings in this area.

1.1 Role of phosphates in biological systems

Phosphates are among the fundamental building blocks that play a central role in life on Earth. They form the basis for both the storage and transfer of genetic information, as well as the flow of metabolic energy within biological systems. The ubiquitous nature of phosphate esters and anhydrides - such as those found in deoxyribonucleic acid (DNA), ribonucleic acid (RNA), adenosine triphosphate (ATP), and polyphosphate (polyP) - highlights their fundamental importance [1]. Some of the phosphates found in biological systems and their respective functions are summarised in Table 1.1.

A key characteristic enabling these roles is the ability of phosphoric acid to link molecular units while retaining an ionisable group. This inherent negative charge at physiological pH serves a dual purpose: it helps to retain these molecules within cellular boundaries defined by lipid membranes, and more importantly, it confers kinetic stability upon phosphate esters and anhydrides by electrostatically repelling nucleophilic attack, particularly from water [1]. For instance, the half-time for hydrolysis at 25 °C for a phosphomonoester monoanion (P-O) is about 90 years; however, for a phosphodiester anion (P-O), this number increases dramatically to approximately 16 million years

[2]. This stability is crucial for maintaining the integrity of genetic material but can be readily overcome by enzymatic catalysis when there is a metabolic demand.

Phosphates are involved in numerous processes in living systems, such as cell signalling and sensation, regulation of metabolism, blood coagulation, and bone formation [3, 4]. Their role is perhaps most evident in cellular energetics, where ATP functions as the universal energy currency. The energy derived from nutrients such as glucose is captured and stored within the high-energy phosphoanhydride bonds linking the phosphate groups of ATP. This energy is released upon hydrolysis of the terminal phosphoanhydride bond (P-O bond between β and γ in Figure 1.1), typically yielding adenosine diphosphate (ADP) and inorganic phosphate (P_i). The cleavage of this bond provides the thermodynamic driving force for the majority of cellular processes, including biosynthesis, active transport, and mechanical work such as muscle contraction. The standard free energy change for ATP hydrolysis is substantial ($\Delta G^0 = -30.5 \text{ kJ mol}^{-1}$), and under cellular conditions, the actual free energy release is often considerably greater. Specifically, the experimentally obtained ΔG values are approximately -59 to $-53.5 \text{ kJ mol}^{-1}$ in the liver and about -61.7 to $-59.5 \text{ kJ mol}^{-1}$ in the heart [3].

Beyond ATP, polyP - a linear polymer of orthophosphate residues linked by similar high-energy phosphoanhydride bonds - represents another significant phosphate-based energy storage found across all domains of life, including mammalian cells. However, in mammalian cells, the concentration of polyP is significantly lower compared to that in microorganisms. While its roles in mammals are still being fully elucidated,

Phosphate	Biological role
DNA/RNA	Genetic material
ADP/ATP	Intracellular energy transfer
cAMP	Cellular signalling
Polyphosphate	Energy storage, Cellular signalling
Creatine phosphate	Intracellular energy transfer
Phosphoenolpyruvate	Metabolism
Pyridoxal phosphate	Coenzyme
Nicotinamide adenine dinucleotide	Calcium signalling
Fructose 1,6-diphosphate	Metabolism
Glucose-6-phosphate	Metabolism
Isopentenyl pyrophosphate	Metabolism
Ribose-6-phosphate	Metabolism
Glycerol 3-phosphate	Metabolism
Dihydroxyacetone phosphate	Calvin cycle, metabolism
Inositol phosphates	Cellular signalling

Table 1.1: Examples of biologically relevant phosphates and their roles. Reproduced and adapted from [5].

dated, polyP metabolism is intrinsically linked to the cellular energy status. Mitochondrial polyP levels fluctuate with respiratory activity and appear to depend on F_0F_1 -ATP synthase function, suggesting a role in mitochondrial bioenergetics, potentially acting as an energy reservoir [6].

The efficient transfer of energy stored in phosphate bonds from sites of production (e.g., mitochondria) to sites of utilisation (e.g., ATPases involved in muscle contraction or ion transport) is crucial. Simple diffusion of ATP is often insufficient due to the complexity of intracellular architecture and the potential for large concentration gradients to arise, which would be thermodynamically inefficient. Instead, cells employ phosphotransfer networks, utilising enzymes such as creatine kinase and adenylate kinase, which catalyse phosphoryl exchange reactions. These networks act as 'phosphoryl wires', facilitating the efficient conduction of high-energy phosphoryl groups and energetic signals throughout the cell with minimal energy dissipation or accumulation of inhibitory products such as ADP. The existence of these networks underscores the dynamic and highly organised nature of cellular energy management, where phosphates - mainly in the form of ATP - serve as the key energy carriers [7].

The synthesis of ATP occurs primarily through oxidative phosphorylation in mitochondria, a process tightly coupled to the electron transport chain, which establishes a proton-motive force (Δp) across the inner mitochondrial membrane. This electrochemical potential energy is used by the molecular machine ATP synthase. Interestingly, the principal energy input required by ATP synthase is not for the chemical formation of the phosphoanhydride bond itself, but rather for the conformational changes necessary to release the newly synthesised, tightly bound ATP molecule from the enzyme's catalytic site. This 'binding change mechanism' involves the cooperative, sequential action of the enzyme's multiple catalytic sites, driven by proton flow. The hydrolysis of ATP to ADP and P_i is catalysed by a variety of enzymes, including ATPases and possibly F_1 -ATPase, which are frequently coupled to other cellular processes [8].

In summary, the unique chemical properties of phosphates - their ability to form stable esters and energy-rich anhydrides, along with their negative charge - combined with the evolution of sophisticated enzymatic machinery for their synthesis, transfer, and hydrolysis, have secured their vital role in virtually all life processes.

1.2 Enzymes involved in phosphate hydrolysis

The hydrolysis of high-energy phosphoanhydride bonds, particularly the terminal bond in ATP, is a cornerstone of cellular bioenergetics. While numerous enzymes utilise ATP hydrolysis, the F_0F_1 -ATP synthase complex, primarily known for synthesising ATP,

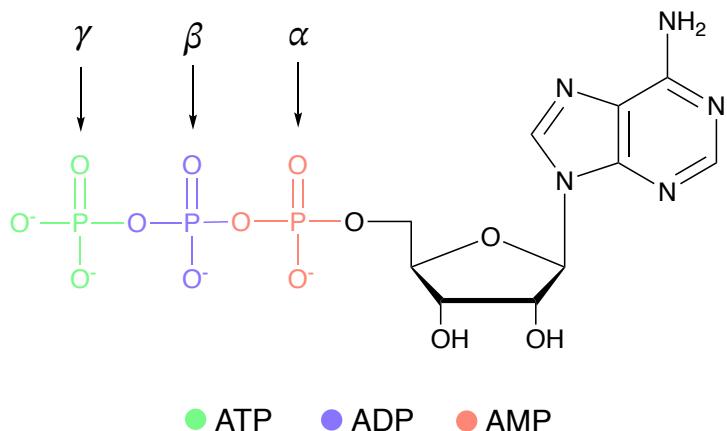


Figure 1.1: Chemical structures of the AMP, ADP, and ATP molecules with the phosphates marked as α , β , and γ , respectively.

also shows potential ATP hydrolytic activity, particularly through its F_1 component (F_1 -ATPase). This enzyme complex, therefore, plays a dual role in managing the cell's primary energy currency [8–10]. Furthermore, recent evidence suggests that this complex may also participate in the metabolism, including the hydrolysis, of polyP in mammalian cells [11, 12].

The F_0F_1 -ATP synthase is a molecular motor embedded in the mitochondrial membrane. It consists of two major domains: the F_1 domain, which carries the catalytic sites, and the F_0 domain, which is embedded within the membrane. These domains are connected by a central rotor stalk and a peripheral stator stalk [10, 13, 14]. The activity of this enzyme is coupled with the electron-transport chain, as illustrated in Figure 1.2.

The F_1 domain ($\alpha_3\beta_3\gamma\delta\epsilon$ stoichiometry) extends into the mitochondrial matrix. It has a globular shape as can be seen in Figure 1.2. The catalytic sites for ATP synthesis and hydrolysis are located on the three β subunits, which interact with the α subunits. When functioning in reverse, the F_1 domain acts as an F_1 -ATPase, hydrolysing ATP. This hydrolysis drives the counterclockwise rotation (as viewed from the membrane) of the central stalk, composed of the γ , δ , and ϵ subunits [8, 10, 13]. If coupled to the F_0 domain, this rotation actively pumps protons from the matrix, thereby generating or maintaining the proton-motive force (Δp). This reverse function is especially important under conditions of low Δp , where it helps prevent its complete dissipation at the expense of cellular ATP and possibly polyP [9, 10, 12].

The mechanism of ATP hydrolysis (cleavage of the P-O bond between β and γ in Figure 1.1) follows the principles of the binding change mechanism [10]. The rotation of the asymmetric γ subunit induces sequential conformational changes in the three β subunits, cycling them through states analogous to those in synthesis: an 'open' state that binds ATP, a 'tight' state that facilitates hydrolysis, and a subsequent 'open' state

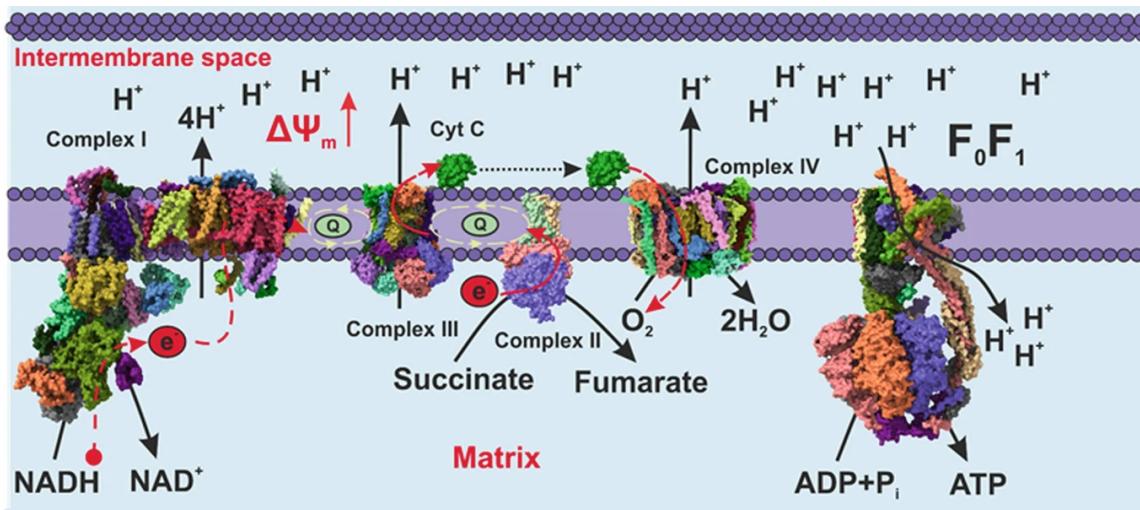


Figure 1.2: Electron transport chain coupled with oxidative phosphorylation in mitochondria. Reproduced from [11].

that releases ADP and P_i [8, 13]. The hydrolysis of each ATP molecule is associated with a 120° rotation of the central stalk, which occurs in substeps [10].

While the metabolism of inorganic polyP is well-characterised in microorganisms via specific kinases (PPK) and phosphatases (PPX), the enzymes responsible for its turnover in mammalian cells remain largely unknown. Recent studies using immuno-captured F₀F₁-ATPase have demonstrated that the enzyme complex can hydrolyse polyP. This polyP hydrolysis appears to drive the enzyme's proton-pumping activity, akin to ATP hydrolysis, and is sensitive to oligomycin, a specific F₀F₁-ATP synthase inhibitor. Medium- and long-chain polyP molecules, made of 60 and 130 orthophosphate units, respectively, seem to be effective substrates for this hydrolytic activity. Docking simulations support the feasibility of polyP binding to the nucleotide-binding sites within the F₁ domain. This suggests that polyP could serve as an alternative energy source for the F₀F₁ complex, potentially helping to maintain mitochondrial membrane potential when ATP levels are compromised [11, 12].

There is growing confidence that the F₁-ATPase could act not only as an ATP hydrolase but also potentially as a polyP hydrolase. However, other enzymes contribute to phosphate metabolism as well. In the context of polyP, mammalian enzymes such as alkaline phosphatase (ALP) have demonstrated exopolyphosphatase activity, capable of degrading polyP chains of various lengths [11].

The world of enzymes - and phosphate hydrolysis by F₁-ATPase in particular - is both fascinating and complex. The F₁-ATPase is a molecular machine capable of hydrolysing ATP and polyP, yet the precise mechanism of hydrolysis remains not well understood. In order to address this gap, it is necessary to investigate the fundamental reaction mechanisms of phosphate hydrolysis, beginning with the simplest phosphate

esters in less complex environments such as bulk water.

1.3 Reaction mechanism: phosphate monoesters

Computational and experimental studies have provided significant insights into the mechanisms of phosphate hydrolysis reactions. Various systems and methodologies have been employed to explore the details of these fundamental biological processes. The debate often centres on whether the reaction proceeds via an associative mechanism (bond formation precedes bond breaking) or a dissociative mechanism (bond breaking precedes bond formation), and the nature of the proton transfer.

1.3.1 Phosphates

Starting from the simplest possible system, it has been shown that the hydrolysis of methyl phosphate dianion ($\text{MeMP}^{\cdot\cdot-}$) in water can proceed via either associative or concerted mechanisms [5, 15–17]. A schematic representation of these mechanisms is presented in Figure 1.3, which illustrates the More O’Ferrall-Jencks (MFJ) diagram. The MFJ plot is a useful two-dimensional graphical representation of multidimensional free energy surfaces.

The associative mechanism may proceed in two ways: stepwise ($\text{A}_N + \text{D}_N$, where A_N stands for nucleophilic addition and D_N for nucleophilic departure) or concerted (A_ND_N). The stepwise mechanism involves two transition states and an intermediate. In contrast, the concerted mechanism proceeds through a single transition state without the formation of intermediates [17].

In the case of the associative/stepwise mechanism ($\text{A}_N + \text{D}_N$, Figure 1.4), the nucleophile (Nuc) approaches the phosphorus atom while the leaving group (lg) is still attached. Upon the nucleophile’s approach, a concerted proton transfer (PT) occurs to one of the non-bridging oxygens. The reaction proceeds through a compact pentaco-ordinated transition state with a trigonal bipyramidal geometry, followed by a compact intermediate and the elimination of the leaving group in a subsequent transition state.

Regarding the associative/concerted mechanism (A_ND_N), it proceeds in a manner quite similar to the first step of the associative/stepwise pathway. The reaction also involves a compact transition state in which bond formation and bond cleavage occur simultaneously.

It has been shown that the protonation state of methyl phosphate lowers the overall barrier height of the rate-limiting step; however, it does not alter the reaction mechanism [18]. For the methyl phosphate monoanion ($\text{MeHMP}^{\cdot+}$), the calculated barrier

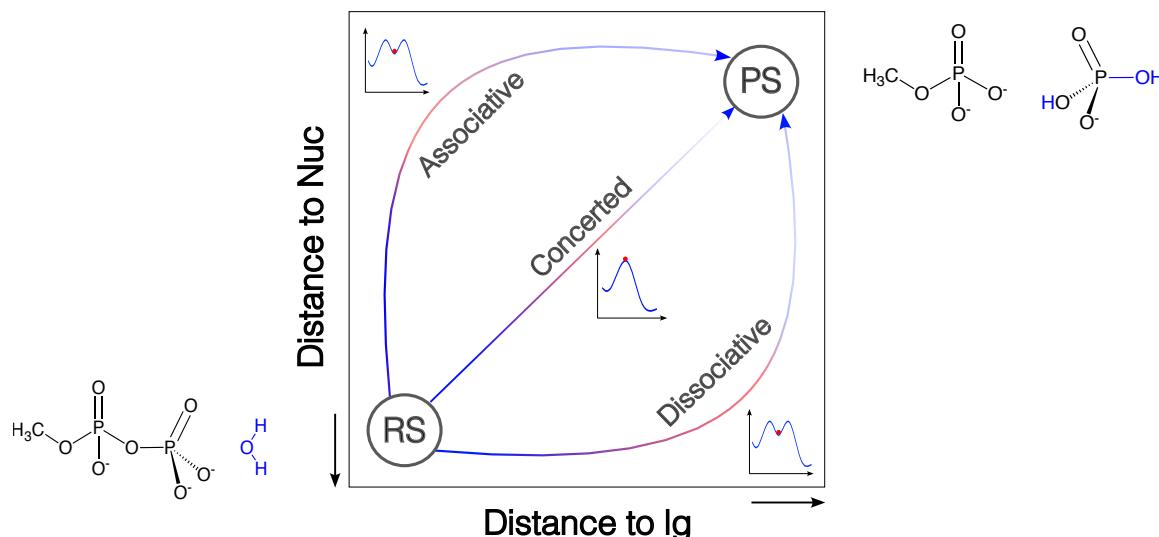


Figure 1.3: More O'Ferrall-Jencks (MFJ) plot of the possible reaction mechanisms for phosphate hydrolysis. The plot shows the free energy as a function of two reaction coordinates: the distance between phosphorus and the nucleophile (Nuc), and the distance between the leaving group (Ig) and the phosphorus atom. RS stands for reactant state, PS for product state.

height $\Delta G_{\text{calc}}^\ddagger$ is approximately 6-7 kcal/mol lower than that of the MeMP. A similar effect was observed when OH^- acted as a nucleophile instead of a water molecule [16] (40 vs 47 kcal/mol, respectively). The latter fact arises a question about the proton-transfer in this reaction, for instance, whether it could happen in a concerted or a step-wise manner, in which the PT happens in the pre-equilibration phase.

For the associative mechanism, calculated barrier heights $\Delta G_{\text{calc}}^\ddagger$ lie in the range of 33.7-47.2 kcal/mol, while experimental values obtained at 25 °C range between 30.6 and 44.3 kcal/mol, depending on the protonation state. Detailed information about the calculated and experimentally determined barrier heights can be found in Table 1.2. Corresponding data on transition state structures and intermediates is summarised in Table 1.3.

The concerted mechanism ($A_N D_N$) is characterised by a single transition state where the nucleophile approaches the phosphorus atom while the leaving group remains attached. The reaction proceeds via a compact pentacoordinated transition state with a trigonal bipyramidal geometry, which is more expansive compared to that of the associative mechanism (Table 1.3). In this transition state, the distance between the phosphorus atom and the nucleophile is approximately 2.06-2.75 Å, while the distance between the leaving group and the phosphorus atom is around 2.61-2.75 Å. The barrier heights for the concerted mechanism are 44.5 and 47 kcal/mol (Table 1.2).

As can be observed, it is rather difficult to clearly distinguish between the associative and concerted mechanisms, and it appears that both may occur in bulk water. Even by looking at the activation entropies of both reaction pathways, it's clear that the

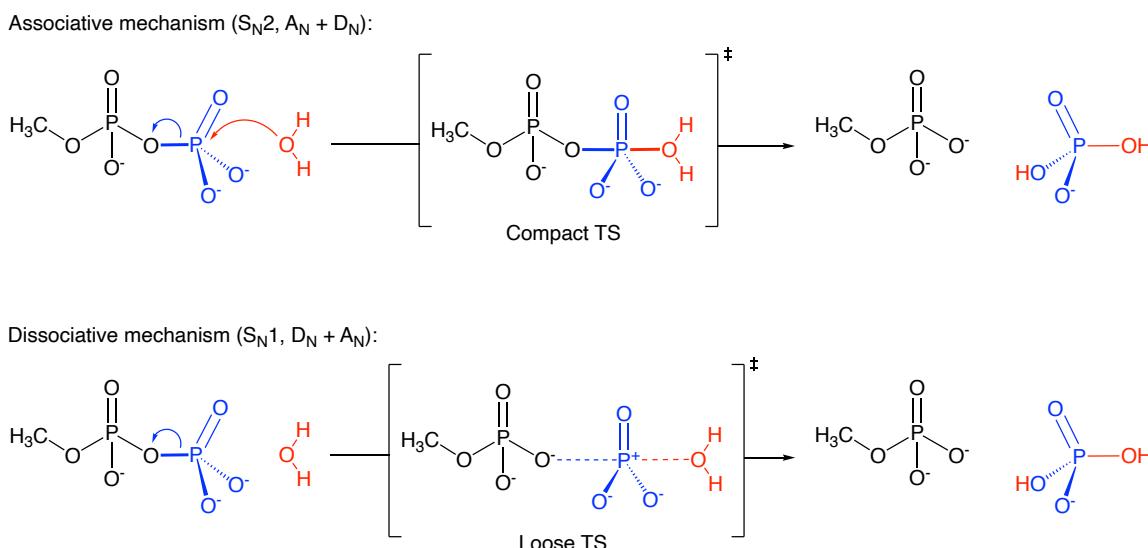


Figure 1.4: Associative and dissociative reaction mechanisms. For the transition states, only one of the two is shown. The nucleophile (Nuc) is shown in red, the leaving group (lg) in black, and the phosphoryl group in yellow.

corresponding values are similarly small: 0.7 and -1.6 kcal/mol for the associative and concerted mechanisms, respectively [17]. Nevertheless, the dissociative mechanism is unlikely to take place, or at least it has not been observed.

1.3.2 Diphosphates

Moving on to more complex systems, the hydrolysis of pyrophosphates (diphosphates), namely methyl diphosphate trianion (MeDP), has also been thoroughly studied. It has been shown that the reaction mechanism can proceed through either associative or concerted pathways, just as in the case of methyl phosphate. In contrast, the associative mechanism has been demonstrated to be solely concerted. However, there is a twist to this story: the dissociative mechanism has also been observed, in both concerted and stepwise forms [15, 16, 19]. A schematic representation of the mechanisms can be found in Figures 1.3 and 1.4.

In the associative/concerted mechanism, the reaction undergoes the same steps as discussed in Subsection 1.3.1. The transition state has a similarly compact geometry: the distance between the phosphorus atom and the nucleophile is approximately 2.03–2.26 Å, and the distance between the leaving group and the phosphorus atom is around 1.82–2.4 Å, as shown in Table 1.3. However, the calculated barrier heights are slightly lower in comparison to methyl phosphate, ranging from 34 to 38 kcal/mol (Table 1.2).

The concerted pathway is characterised by the same general mechanism as discussed in Subsection 1.3.1. The transition state is more expansive than in the associative mechanism, with the distance between the phosphorus atom and the nucleophile

Table 1.2: Summary of computational and experimental studies on phosphate hydrolysis. In the case of calculated ΔG^\ddagger , the rate-limiting step is given. ¹The values were calculated using the transition state theory (TST).

System	Method	Level of theory	Mechanism	ΔG^\ddagger (kcal/mol)	Ref.
MeMP ²⁻ + H ₂ O	DFT	B3LYP/6-311+G** and COSMO	Associative Concerted	47.2 44.5	[15]
MeMP ²⁻ + H ₂ O	DFT	B3LYP/6-311++G** and COSMO	Associative	47	[16]
MeMP ²⁻ + 3 H ₂ O	DFT	M06-2X/6-311+G** and SMD	Associative Concerted	\approx 36 \approx 44	[17]
MeMP ²⁻ + 4 H ₂ O	DFT	M06-2X/6-311+G**	Associative	\approx 40.8 \pm 1.9	[18]
MeHMP ⁻ + 4 H ₂ O	DFT	M06-2X/6-311+G**	Associative	\approx 33.7 \pm 1.7	[18]
MeDP ³⁻ + 2 H ₂ O	DFT	B3LYP/6-311++G** and PCM	Associative Dissociative	34.64 35.24	[19]
MeDP ³⁻ + H ₂ O	DFT	B3LYP/6-311+G** and COSMO	Associative Dissociative	34.8 30.3	[15]
MeDP ³⁻ + H ₂ O	DFT	B3LYP/6-311++G** and COSMO	Associative Concerted	38 34	[16]
MeHDP ²⁻ + H ₂ O	DFT	B3LYP/6-311++G** and COSMO	Associative Concerted	34 31	[16]
MeDP ³⁻ + Mg ²⁺ + 5 H ₂ O	QM/MM, FEP (EVB)	B3LYP/6-311++G** and MM	Associative Concerted Dissociative	35 34 35	[16]
MeTP ⁴⁻ + Mg ²⁺ + 54 H ₂ O	CPMD	PBE/PW with Troullier-Martins pseudopotentials	Associative Concerted Dissociative	39.1 35.1 36.6	[20]
MeTP ⁴⁻ + Mg ²⁺ + 113 H ₂ O	BOMD, metadynamics	BLYP/TZV2P with GTH pseudopotentials	Associative Concerted	29 29-30	[21]
ATP ⁴⁻ + Mg ²⁺ + 4163 H ₂ O + counterions	QM/MM, NEB	B3LYP/6-311++G** and MM	Concerted	32.5	[22]
ATP ⁴⁻ + Mg ²⁺ + 1800 H ₂ O + counterions	QM/MM, QM = CPMD	BLYP/PW with Troullier-Martins pseudopotentials and MM	Associative Dissociative	36.2 33.4	[23]
Methyl phosphate dianion	Exp. at 25°C	–	–	44.3	[2]
Methyl phosphate monoanion	Exp. at 25°C	–	–	30.6	[2]
Pyrophosphate trianion	Exp. at 25°C	–	–	29.2	[2]
Pyrophosphate dianion	Exp. at 25°C	–	–	27.7	[2]
ADP ²⁻ (or ATP ³⁻)	Exp. at 25°C	–	–	27.5	[2]
ATPH ³⁻ (or ATP ⁴⁻)	Exp. at 70°C	pH=6.69-7.66	–	24.34-24.78 ¹	[24]
dADPH ²⁻ (or dADP ³⁻)	Exp. at 70°C	pH=6.82	–	24.25 ¹	[25]
dATPH ³⁻ (or dATP ⁴⁻)	Exp. at 70°C	pH=7.00	–	24.50 ¹	[25]
MgATPH ⁻ (or MgATP ²⁻)	Exp. at 70°C	pH=6.59-7.63	–	24.59-24.64 ¹	[24]
CaATPH ⁻ (or CaATP ²⁻)	Exp. at 70°C	pH=6.67-7.01	–	25.71-25.72 ¹	[24]

being approximately 2.26–2.5 Å, while the distance between the leaving group and the phosphorus atom is around 2.7 Å (Table 1.3). The barrier heights for the concerted mechanism are 31 and 34 kcal/mol (Table 1.2), which is notably lower than in the case of methyl phosphate.

In general, it can be noted that the barrier height is strongly dependent on the pK_a value of the leaving group. The lower the pK_a , the lower the ΔG^\ddagger ($pK_a(\text{CH}_3\text{O}^-) = 15.5$ vs $pK_a(\text{CH}_3\text{PO}_4^{2-}) = 6.3$). Not only does a lower pK_a reduce the barrier height, but it also favours the concerted and dissociative mechanisms [16].

The dissociative mechanism can proceed via both concerted and stepwise routes. The dissociative/concerted pathway is quite similar to the general concerted mechanism. The main difference lies in the synchrony of the transition state: while the general concerted mechanism has a more synchronous transition state, the dissociative/concerted one features a greater distance between the phosphorus atom and the leaving group.

On the other hand, the dissociative/stepwise pathway ($D_N + A_N$) is characterised by the departure of the leaving group from the phosphorus atom before the nucleophile approaches. Thus, there is clearly no bond remaining between the phosphorus and the leaving group. Following the departure of the leaving group, a planar metaphosphate PO_3^- is formed, as illustrated in Figure 1.4. The transition state is more expansive compared to that of the associative mechanism, with the distance between the phosphorus atom and the nucleophile being 2.7 Å, and the distance between the leaving group and the phosphorus atom being 3.4 Å (Table 1.3). Consequently, after TS_1 , the system reaches an intermediate step in which the nucleophile is positioned closer to the metaphosphate, followed by an attack and bond formation in TS_2 .

The barrier heights for the dissociative mechanism lie in the range of 30.3–35.24 kcal/mol (Table 1.2), which is lower than those for any of the previously mentioned mechanisms. Comparing the calculated and experimentally obtained ΔG^\ddagger clearly indicates that the dissociative mechanism is more favourable than the associative and concerted ones. The $\Delta G_{\text{exp}}^\ddagger$ values for the pyrophosphate trianion and dianion are 29.2 and 27.7 kcal/mol, respectively. The influence of one-water (1W) or two-water (2W) mechanisms has also been explored [19], but the overall barriers remain similar.

The dissociative mechanism is further favoured by the presence of metal ions, such as Mg^{2+} , as well as in cases where MeDP is protonated, i.e. methyl diphosphate dianion (MeHDP). Interestingly, in the latter case, the proton always transfers to the leaving group en route to the product state [16].

Even though computational studies suggest that Mg^{2+} promotes the dissociative mechanism, experimental data do not support this hypothesis [24], since the $\Delta G_{\text{exp}}^\ddagger$ obtained at 70 °C shows little to no difference, at least in the case of adenosine triphos-

Table 1.3: Summary of the distances between the phosphorus atom and the nucleophile as well as the leaving group in the transition states and intermediates. All distances are in Å.

System	Mechanism	TS ₁		Intermediate		TS ₂		Ref.
		d(P-O _{Nuc})	d(P-O _{lg})	d(P-O _{Nuc})	d(P-O _{lg})	d(P-O _{Nuc})	d(P-O _{lg})	
MeMP ²⁻	Associative (A _N D _N)	2.0	1.8	—	—	—	—	[16]
	Associative (A _N D _N)	1.9	2.15	—	—	—	—	[15]
	Associative (A _N + D _N)	2.16	1.71	1.84	1.78	1.71	2.24	[17]
	Associative (A _N + D _N)	2.08	1.78	1.99	1.80	1.77	2.52	[18]
	Concerted (A _N D _N)	2.75	2.75	—	—	—	—	[15]
	Concerted (A _N D _N)	2.06	2.61	—	—	—	—	[17]
	Associative (A _N + D _N)	2.26	1.66	1.76	1.77	1.68	2.25	[18]
MeHMP ⁻	Associative (A _N D _N)	2.2	2.0	—	—	—	—	[15]
	Associative (A _N D _N)	2.03	1.88	—	—	—	—	[16]
	Associative (A _N D _N)	2.2	2.0	—	—	—	—	[19]
	Concerted (A _N D _N)	2.5	2.7	—	—	—	—	[16]
	Dissociative (A _N D _N)	2.8	3.25	—	—	—	—	[15]
	Dissociative (D _N + A _N)	2.7	3.4	2.0	3.75	1.7	3.75	[19]
	Associative (A _N D _N)	2.1	2.4	—	—	—	—	[16]
MeDP ³⁻ + Mg ²⁺	Concerted (A _N D _N)	2.3	2.7	—	—	—	—	[16]
	Dissociative (A _N D _N)	2.8	3.4	—	—	—	—	[16]
	Associative (A _N D _N)	2.26	1.82	—	—	—	—	[16]
MeHDP ²⁻	Concerted (A _N D _N)	2.26	2.78	—	—	—	—	[16]
	Associative (A _N D _N)	1.9	2.0	—	—	—	—	[20]
MeTP ⁴⁻ + Mg ²⁺	Associative (A _N D _N)	2.03	3.11	1.95	3.06	1.66	3.26	[21]
	Associative (A _N + D _N)	2.5	2.6	—	—	—	—	[20]
	Concerted (A _N D _N)	2.28	2.69	—	—	—	—	[21]
	Dissociative (A _N D _N)	3.6	3.5	—	—	—	—	[20]
	Associative (A _N D _N)	1.9	1.9	—	—	—	—	[23]
ATP ⁴⁻ + Mg ²⁺	Concerted (A _N D _N)	2.8	3.2	—	—	—	—	[22]
	Dissociative (A _N D _N)	3.5	3.5	—	—	—	—	[22]

phate (Table 1.2).

1.3.3 Triphosphates

Last but not least, let us consider the hydrolysis of triphosphates. These systems more closely resemble the biological environment, particularly the processes associated with energy metabolism.

The hydrolysis of methyl triphosphate tetranion (MeTP) and ATP has been studied using a range of computational methods. It has been shown that the reaction mechanisms share many similarities with those observed in mono- and diphosphates. Specifically, the mechanism may proceed via associative/concerted and associative/stepwise routes, as well as concerted and dissociative/concerted pathways (Table 1.3).

When comparing the calculated and experimentally obtained ΔG^\ddagger values, it is difficult to clearly distinguish between the aforementioned mechanisms. The calculated barrier heights span a range from 29 to 39.1 kcal/mol, whereas experimental data suggest that the barrier height for ATP should be around 27.5 kcal/mol, as shown in Table 1.2. The more complex the system becomes, the more factors one must likely take into account.

In summary, computational investigations reveal a nuanced and peculiar picture of phosphate hydrolysis. The preferred mechanism (associative, concerted, or dissociative) and the proton transfer route (1W, 2W, etc.) depend significantly on specific factors such as the pKa of the leaving group, the protonation state, the presence of metal ions like Mg^{2+} , and the surrounding solvent environment.

Moreover, it is important to keep in mind that in order to properly study the underlying free energy surface, adequate sampling of the reaction space is crucial. To achieve this, various free energy techniques such as metadynamics can be employed, provided that the level of theory is sufficient to describe a system of realistic size while allowing the results to be obtained within a reasonable timeframe. This is precisely the goal of the present project.

1.4 Research goals

To begin with, I'd like to quote the following line from Paul A. M. Dirac [26], which I find particularly relevant to the topic of this thesis:

The underlying physical laws necessary for the mathematical theory of a large part of physics and the whole of chemistry are thus completely known, and the difficulty is only that the exact application of these laws leads to

equations much too complicated to be soluble. *It therefore becomes desirable that approximate practical methods of applying quantum mechanics should be developed, which can lead to an explanation of the main features of complex atomic systems without too much computation.*

Chapter 2

Theory

2.1 A brief introduction to statistical mechanics

The discussion in this section is mostly based on the “Introduction to Computational Chemistry” textbook written by Jensen [27], “Statistical Mechanics: Theory and Molecular Simulation” by Tuckermann [28], and “Understanding Molecular Simulation: From Algorithms to Applications” by Frenkel and Smit [29] unless stated otherwise.

2.1.1 Partition functions

The development of the field of statistical mechanics has been crucial for the computational chemistry community, as it enables the connection between the jigglings and wiggles of atoms and the properties of much larger systems such as liquids and solids.

Let us begin with the most fundamental concept: the partition function. The partition function is akin to a Swiss army knife in statistical mechanics, meaning it is a versatile tool that makes the connection between microscopic and macroscopic properties in thermodynamics possible. In the simplest case of a single molecule, the partition function q takes the following form:

$$q = \sum_{i=\text{levels}}^{\infty} g_i e^{-\epsilon_i/kT} \quad (2.1)$$

Here, it is expressed as a sum over all energy levels ϵ_i of a molecule (or particle), multiplied by a degeneracy factor g_i in cases where multiple levels have the same energy. The term kT represents the Boltzmann factor.

Moving on to a more complex scenario in which the partition function describes multiple molecules, we arrive at the partition function Q for non-interacting particles,

such as those in an ideal gas:

$$Q = q^N \text{ (different particles)} \quad Q = \frac{q^N}{N!} \text{ (identical particles)} \quad (2.2)$$

Here, N denotes the total number of particles. However, one could argue that if we wish to describe a real system such as bulk water, we must account for interactions between molecules. Consequently, Equation 2.2 must be rewritten:

$$Q = \sum_i^{\infty} e^{-E_i/kT} \quad (2.3)$$

In this case, the partition function Q includes contributions from all possible energy states E_i of the system.

Although the concept of the partition function might initially appear abstract, it can be clarified by expressing it in a different form, namely, within the context of the Rigid-Rotor Harmonic-Oscillator (RRHO) approximation, where the electronic, vibrational, and rotational degrees of freedom can be separated. For a single molecule case it would look like:

$$q_{\text{tot}} = q_{\text{trans}} \times q_{\text{rot}} \times q_{\text{vib}} \times q_{\text{elec}} \quad (2.4)$$

Let us now examine each contribution in more detail. From this point onward we will consider polyatomic molecules in the formulation of the partition functions, unless stated otherwise.

The translational partition function q_{trans} can be derived from the energy expression for a particle in a one-dimensional box and is given by:

$$q_{\text{trans}} = \left(\frac{2\pi MkT}{h^2} \right)^{3/2} V \quad (2.5)$$

Here, M is the total molecular mass, and V is the volume. Turning to the rotational partition function q_{rot} , it can be derived from the Schrödinger equation for a diatomic "rigid rotor" and has the following form:

$$q_{\text{rot}} = \frac{8\pi^2 I k T}{h^2 \sigma} \quad (2.6)$$

In this expression, I denotes the moment of inertia, and σ represents the symmetry factor, i.e. the order of the rotational subgroup within the molecular point group. For polyatomic molecules, writing an exact expression is more complex, but an approximate form can be used:

$$q_{\text{rot}} = \frac{\sqrt{\pi}}{\sigma} \left(\frac{8\pi^2 kT}{h^2} \right)^{3/2} \sqrt{l_1 l_2 l_3} \quad (2.7)$$

For the vibrational partition function q_{vib} , it is expressed as a product over the various vibrational modes of a molecule, each with frequency ν_i :

$$q_{\text{vib}} = \prod_i \frac{e^{-h\nu_i/2kT}}{1 - e^{-h\nu_i/kT}} \quad (2.8)$$

Lastly, the electronic partition function q_{elec} is given as a sum over all electronic states of a molecule, from the ground state to all excited states. However, since the energy difference between the ground state and higher states is usually much greater than kT at ambient temperatures, the function can typically be approximated by considering only the ground state:

$$q_{\text{elec}} = \sum_{i=0}^{\infty} g_i e^{-\epsilon_i/kT} \approx g_0 e^{-\epsilon_0/kT} \quad (2.9)$$

This contribution can thus be calculated by solving the electronic Schrödinger equation.

2.1.2 Macroscopic properties and thermodynamic functions

Once the partition function is determined, it provides a direct means of evaluating macroscopic properties. For instance, the internal energy U and the Helmholtz free energy A can be calculated from the partition function Q :

$$U = kT^2 \left(\frac{\partial \ln Q}{\partial T} \right)_V \quad (2.10)$$

$$A = -kT \ln Q \quad (2.11)$$

In addition, other macroscopic properties, such as pressure P and the heat capacity at constant volume C_V , can also be expressed in terms of the partition function:

$$P = - \left(\frac{\partial A}{\partial V} \right)_T = kT \left(\frac{\partial \ln Q}{\partial V} \right)_T \quad (2.12)$$

$$C_V = \left(\frac{\partial U}{\partial T} \right)_V = 2kT \left(\frac{\partial \ln Q}{\partial T} \right)_V + kT^2 \left(\frac{\partial^2 \ln Q}{\partial T^2} \right)_V \quad (2.13)$$

Turning to thermodynamic functions, namely enthalpy H , entropy S , and Gibbs free

energy G , these can also be derived from the partition function Q :

$$H = U + PV = kT^2 \left(\frac{\partial \ln Q}{\partial T} \right)_V + kTV \left(\frac{\partial \ln Q}{\partial V} \right)_T \quad (2.14)$$

$$S = \frac{U - A}{T} = kT \left(\frac{\partial \ln Q}{\partial T} \right)_V + k \ln Q \quad (2.15)$$

$$G = H - TS = kTV \left(\frac{\partial \ln Q}{\partial V} \right)_T - kT \ln Q \quad (2.16)$$

This connection between macroscopic observables, thermodynamic functions, and the partition function once again highlights its fundamental importance in statistical thermodynamics.

2.1.3 The canonical ensemble

Having established a method to calculate the macroscopic properties of a system we implicitly relied on averaging over a large enough number of states. Therefore one may naturally ask: how can we sample enough configurations to apply the equations described in the previous section under conditions that resemble those in experiments? One such answer is the canonical ensemble.

The canonical ensemble describes a system at constant temperature T , fixed volume V , and a fixed number of particles N (NVT). In this ensemble, the system is in contact with a heat bath, which makes it particularly relevant to most molecular simulations that are describing the experimental conditions, where the temperature is externally controlled while the internal energy of the system is allowed to fluctuate.

Since the energy fluctuates in the canonical ensemble, a logical step is to estimate the magnitude of these fluctuations:

$$\frac{\Delta E}{E} \sim \frac{\sqrt{N}}{N} \sim \frac{1}{\sqrt{N}} \quad (2.17)$$

Here, N denotes the number of particles, and thus for sufficiently large systems, the relative energy fluctuations become negligible.

The use of the canonical ensemble implicitly assumes that the system is ergodic, meaning that time averages obtained from simulation trajectories are equivalent to ensemble averages over the Boltzmann distribution. This assumption is, for instance, central to molecular dynamics simulations where the canonical ensemble can be sampled.

2.1.4 Classical forcefields and molecular dynamics

2.1.5 Enhanced sampling techniques

$$V_G(S(x), t) = w \sum_{t'=\tau_G, 2\tau_G, \dots}^{t' < t} \exp \left(-\frac{(S(x) - s(t'))^2}{2\delta s^2} \right) \quad (2.18)$$

$$\lim_{t \rightarrow \infty} V_G(s, t) \sim -F(s) \quad (2.19)$$

$$V(s, t) = \Delta T \ln \left(1 + \frac{\omega N(s, t)}{\Delta T} \right) \quad (2.20)$$

$$\dot{V}(s, t) = \frac{\omega \Delta T \delta_{s, s(t)}}{\Delta T + \omega N(s, t)} = \omega e^{-[V(s, t)/\Delta T]} \delta_{s, s(t)} \quad (2.21)$$

$$w = \omega e^{-[V(s, t)/\Delta T]} \tau_G \quad (2.22)$$

$$\tilde{F}(s, t) = -\frac{T + \Delta T}{\Delta T} V(s, t) = -(T + \Delta T) \ln \left(1 + \frac{\omega N(s, t)}{\Delta T} \right) \quad (2.23)$$

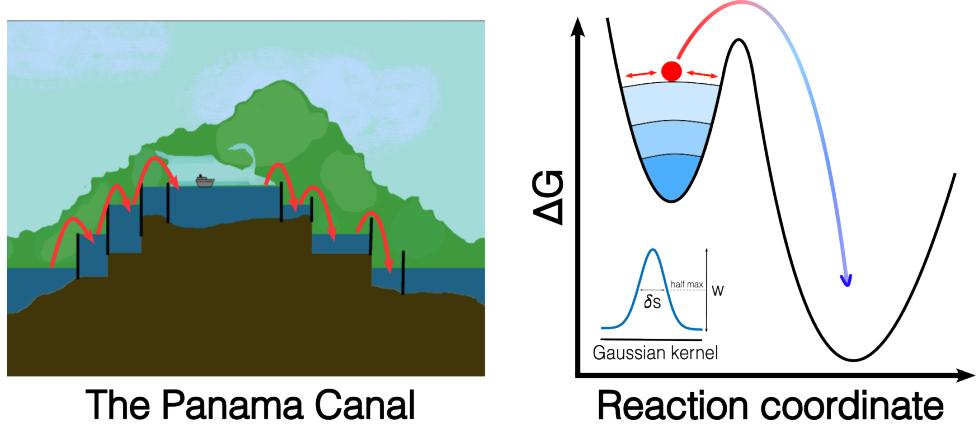


Figure 2.1: Metadynamics. The Panama Canal cartoon reproduced from [30].

2.2 Transition state theory

2.3 Density functional theory

2.3.1 The Kohn-Sham approach

2.3.2 Generalised gradient approximation and PBE functional

2.3.3 *Ab initio* molecular dynamics and GPW method

2.4 Extended tight binding

2.5 Neural network potentials

2.5.1 Message passing graph neural networks

2.5.2 Invariance and equivariance

2.5.3 Equivariant neural network potentials

Chapter 3

Computational details

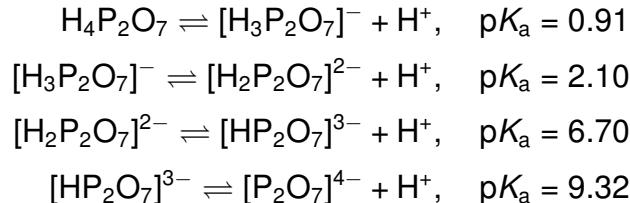
This chapter provides detailed information on the computational methods employed in this work. The first section outlines the generation of the training dataset, including system preparation, initial equilibration using molecular mechanics, exploration of the configuration space at the GFN1-xTB level, further data labelling, and iterative training of the neural network potential. The second section discusses production runs at various temperatures using the fitted neural network potential. The third section describes the workflow for validating the transition states obtained from the simulations, based on the partial Hessian formalism. Finally, the fourth section presents the data analysis and visualisation techniques used to interpret the results.

3.1 Training dataset generation

3.1.1 System preparation

The systems were prepared using the functionality of the CHARMM-GUI webserver [31], specifically the Multicomponent Assembler interface [32].

As a first step, methyl diphosphate trianion (MeDP) and methyl diphosphate dianion (MeHDP) were parameterised using CGenFF [33], i.e., the CHARMM General Force Field. These protonation states were chosen based on the dissociation constants of pyrophosphoric (diphosphoric) acid [34]:



Thus, at physiological pH (7.4), this acid exists in equilibrium between the singly and doubly deprotonated forms. Assuming the methyl group behaves similarly to a proton, the methyl diphosphate molecule was considered to exist as a mixture of the MeHDP and MeDP forms under physiological conditions.

Following successful parameterisation, the system was solvated in a cubic box of TIP3P water molecules, with sodium counterions (Na^+) added to neutralise the system's overall charge. The final system composition is provided in Table 3.1.

3.1.2 Initial equilibration using classical force fields

The system equilibration followed the standard protocol generated by the CHARMM-GUI webserver [31]. Initially, energy minimisation was conducted using the steepest descent algorithm for 5,000 steps.

This was followed by equilibration in the constant number of particles, volume and temperature (NVT) ensemble for 5 ns. During both the minimisation and NVT phases, the solute's heavy atoms were restrained using a harmonic potential with a force constant of $400 \text{ kJ mol}^{-1} \text{ nm}^{-2}$.

Subsequently, the system was equilibrated in the constant number of particles, pressure and temperature (NPT) ensemble for 45 ns. Throughout this procedure, the temperature and pressure were maintained at 300 K and 1 bar, respectively. Temperature was controlled using a ν -rescale thermostat [35] with a coupling constant of 1 ps, and pressure was regulated using an isotropic c -rescale barostat [36] with a coupling constant of 5 ps. A 0.6 nm cut-off was applied for non-bonded interactions, and long-range electrostatics were treated using the Particle Mesh Ewald (PME) method. Periodic boundary conditions (PBC) were applied in all directions throughout the simulation.

All simulations were carried out using GROMACS 2021.4 [37] with CHARMM36m force field [38]. The leap-frog integrator was employed with a time step of 1 fs. All hydrogen-involving bonds were constrained using the LINCS algorithm. The equilibrated box dimensions used for subsequent simulations were taken from the output of the NPT run and are summarised in Table 3.1. Unless otherwise stated, the last frame of the NPT simulations was used as the starting point for all further calculations.

Table 3.1: System composition and simulation box details.

System	Final box dimensions (\AA^3)	No. of H_2O	No. of Na^+	No. of atoms
MeDP	$15.877 \times 15.877 \times 15.877$	119	3	373
MeHDP	$15.901 \times 15.901 \times 15.901$	124	2	388

3.1.3 Collective variables

To effectively sample the reaction space, two types of collective variables (CVs) were employed to bias the system: distances and coordination numbers (CNs). The CN is defined by the following smooth function:

$$\sum_{i \in A} \sum_{j \in B} CN_{ij} = \frac{1 - \left(\frac{r_{ij} - d_0}{r_0} \right)^n}{1 - \left(\frac{r_{ij} - d_0}{r_0} \right)^m} \quad (3.1)$$

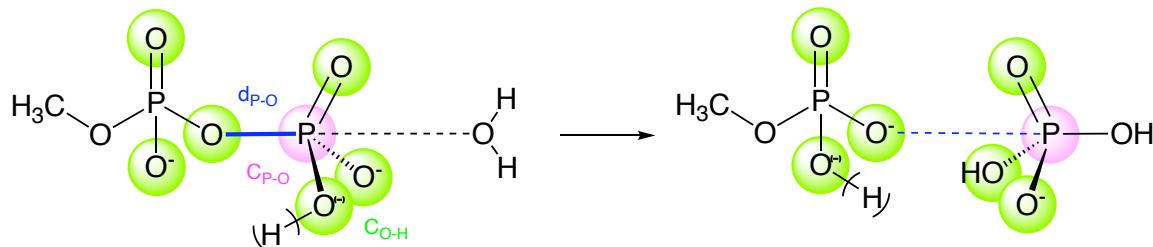
where r_{ij} is the distance between atoms i and j from groups A and B , d_0 is the distance at which the CN begins to decay, r_0 is a characteristic decay length, and n and m are integers that control the steepness of the decay. Typically, $m > n$, ensuring a smooth transition of CN_{ij} from approximately 1 to 0 as the distance increases.

The specific CVs used in this work are shown in Figure 3.1, and their corresponding parameters are as follows:

- Distance between the β -phosphorus and the bridging oxygen (d_{P-O}),
- Coordination number of all water oxygen atoms as well as the bridging oxygen surrounding the β -phosphorus (C_{P-O}): $d_0 = 0$, $r_0 = 2.1 \text{ \AA}$, $n = 8$, $m = 16$,
- Coordination number of non-methyl hydrogen atoms around 5 non-bridging and 1 bridging oxygen atoms (C_{O-H}): $d_0 = 0$, $r_0 = 1.3 \text{ \AA}$, $n = 8$, $m = 16$.

Additionally, the following CVs were monitored to check whether the system was in a reasonable region of the potential energy surface, e.g. to ensure that there is no oxygen exchange between the methyl diphosphate and the water molecules:

- Coordination number of 3 nonbridging oxygens surrounding the β -phosphorus ($C_{P-O_{\text{nonbridging}}}$): $d_0 = 0$, $r_0 = 2.1 \text{ \AA}$, $n = 8$, $m = 16$,



$$d_{P-O} = d(P - O_{lg}), C_{P-O} = CN(P - O_{all}), C_{O-H} = CN(O - H_{all})$$

Figure 3.1: The definition of the collective variables (CVs) used in this work. CN stands for coordination number.

- Coordination number of water oxygens around the β -phosphorus ($C_{P-O_{water}}$): $d_0 = 0$, $r_0 = 2.1 \text{ \AA}$, $n = 8$, $m = 16$,

To avoid sampling of unphysical regions of the potential energy surface, quadratic (harmonic-like) wall potentials were applied to softly constrain certain degrees of freedom. The mathematical form of these wall potentials is given below:

$$\text{For upper walls: } \sum_i k_i \left(\frac{CV_i - a_i + o_i}{s_i} \right)^{e_i} \quad (3.2)$$

$$\text{For lower walls: } \sum_i k_i \left| \frac{CV_i - a_i - o_i}{s_i} \right|^{e_i} \quad (3.3)$$

Here, CV_i denotes the value of the collective variable, k_i is the force constant defining the wall's strength, a_i is the central wall position, o_i is an offset, s_i is a scaling factor, and e_i is the exponent that controls the wall's steepness. When $e_i = 2$, the potential acts harmonically.

The wall potentials applied to the CVs during the simulations are summarised in Table 3.2. The parameters for the wall potentials were chosen based on the expected ranges of the CVs. The force constants were set to ensure that the walls were sufficiently strong to prevent unphysical configurations while allowing for reasonable exploration of the configuration space.

All CV-related computations were performed using the built-in tools of CP2K 2023.1 [39] or PLUMED 2.9.3 [40]. It is important to note that the number and type of CVs, as well as the applied restraints, varied depending on the specific stage of the workflow. In the following sections, the relevant collective variables and wall potentials will be specified accordingly.

Table 3.2: The restraints applied to the collective variables during some of the simulations. In all cases, $o = 0$, $s = 1$, $e = 2$. ¹During the iterative training / production runs. ²Different values for the walls were used depending on the system MeDP/MeHDP. Distances are in \AA and coordination numbers are unitless.

CV	Lower wall	Upper wall	Force constant (kcal mol ⁻¹ Å ⁻²)
d_{P-O}	–	5.0 / 6.0 ¹	500
C_{O-H}	– / 1.3 ²	1.3 / 2.5 ²	1000
$C_{P-O_{\text{nonbridging}}}$	2.6	–	2000
$C_{P-O_{\text{water}}}$	–	1.3	2000

3.1.4 GFN1-xTB based exploration of the configuration space

To generate the initial set of configurations for the training dataset, the system was subjected to molecular dynamics simulations using the semi-empirical GFN1-xTB [41] level of theory. GFN1-xTB provides a good first approximation of the potential energy surface and is computationally efficient, thus making it suitable for relatively long MD simulations of large systems.

Each system was first equilibrated for 5 ps in the NVT ensemble at 300 K to allow the structures to relax at the GFN1-xTB level, including a Grimme's D3 dispersion correction [42]. Following equilibration, we performed 50 ps of well-tempered metadynamics (WTMetaD) [43] simulations in the NVT ensemble. In these simulations, a biasing potential was applied to encourage the system to explore regions of the configuration space beyond the reactant basin. This bias was introduced along two CVs: the distance between the β -phosphorus and the oxygen atom connecting it to the rest of the molecule (d_{P-O}), and the coordination number of all water oxygens together with the bridging oxygen surrounding the β -phosphorus atom (C_{P-O}). No restraints were applied to the system during this stage.

All calculations were carried out using the CP2K 2023.1 package [39] on CPUs. Temperature control was achieved using the ν -rescale thermostat [35], with a time constant of 50 fs during equilibration and 100 fs during the WTMetaD simulations. The self-consistent field (SCF) convergence threshold was set to 10^{-5} a.u. The biasing potential was updated every 25 fs, with a Gaussian hill height of 2 kcal mol⁻¹ and a width of 0.07 for each CV. The bias factor was set to 30. Finally, the integration time step was set to 0.5 fs. Throughout the simulations, PBC were applied in all directions.

3.1.5 Data labeling

All data points were labeled by performing single-point calculations to obtain the energy and force values. These single-point calculations were carried out using the Perdew-Burke-Ernzerhof exchange-correlation functional (PBE) [44], along with a Grimme's D3 dispersion correction and the Becke-Johnson damping function [42, 45]. In all calculations, the Goedecker-Teter-Hutter pseudopotentials (GTH) [46, 47] were used to represent the core electrons, in combination with the triple- ζ valence basis set with two polarisation functions (TZV2P).

The single-point calculations were performed using the Gaussian plane wave method (GPW) implemented in the QUICKSTEP module [48] of the CP2K 2023.1 package [39]. The SCF convergence threshold was set to 10^{-6} a.u. A plane-wave cutoff of 800 Ry was applied for the total density, while a cutoff of 60 Ry was used for

the Kohn-Sham orbitals.

The aforementioned cutoffs were determined based on a convergence test performed on one of the configurations, as described in [49]. An error in total energy of less than 10^{-8} a.u. was considered acceptable for the convergence test. The test was conducted by varying the cutoff for the total density from 400 to 1500 Ry, and the cutoff for the Kohn-Sham orbitals from 10 to 200 Ry. The results of the convergence test are shown in Table A.1.

3.1.6 Iterative training of the neural network potential

We trained a neural network potential (NNP) using the NequIP framework [50], which implements equivariant message-passing networks for atomistic simulations. Regarding the hyperparameters, a radial cutoff distance of 5.0 Å was chosen to describe the atomic environment of the system.

The equivariant part of the neural network was composed of four interaction layers with a maximum tensor rank of $\ell = 1$ or 2 . Feature parity was enabled to include both even and odd components, and 32 features per irreducible representation were used throughout. Scalar and gating nonlinearities were set to `silu` and `tanh` for even and odd parities, respectively. Eight radial basis functions were employed, in combination with a trainable Bessel basis and a polynomial cutoff of order 6.

The invariant subnetwork for radial interaction modelling consisted of two layers with 64 hidden neurons. Self-connections were enabled, and the average number of neighbours was computed automatically based on the dataset.

Training was performed using the Adam optimizer with the AMSGrad variant enabled, and with $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 10^{-8}$. A starting learning rate of 0.01 was used, and the learning rate was adaptively reduced by a factor of 0.5 upon stagnation of the validation loss (patience = 100 epochs). Early stopping was triggered if the validation loss remained unimproved for 50 epochs, if the loss dropped below 1×10^{-5} , or if it exceeded 1×10^4 . The batch size was set to 5. The training was carried out over a period of three days on a single NVIDIA A100 GPU using float64 precision.

To thoroughly sample the reaction space, the training was performed in an iterative manner, where the model was first trained on a small set of data and then used to generate additional data points. This process was repeated until the model converged, with the root mean square error (RMSE) of the atomic forces being less than 40 meV/Å. The workflow is shown in Figure 3.2.

In the end, the full dataset consisted of 12,000 configurations for training and validation, and 1,800 configurations for testing, for both systems (MeDP and MeHDP) combined. This dataset was obtained within the three rounds of iterative training. In

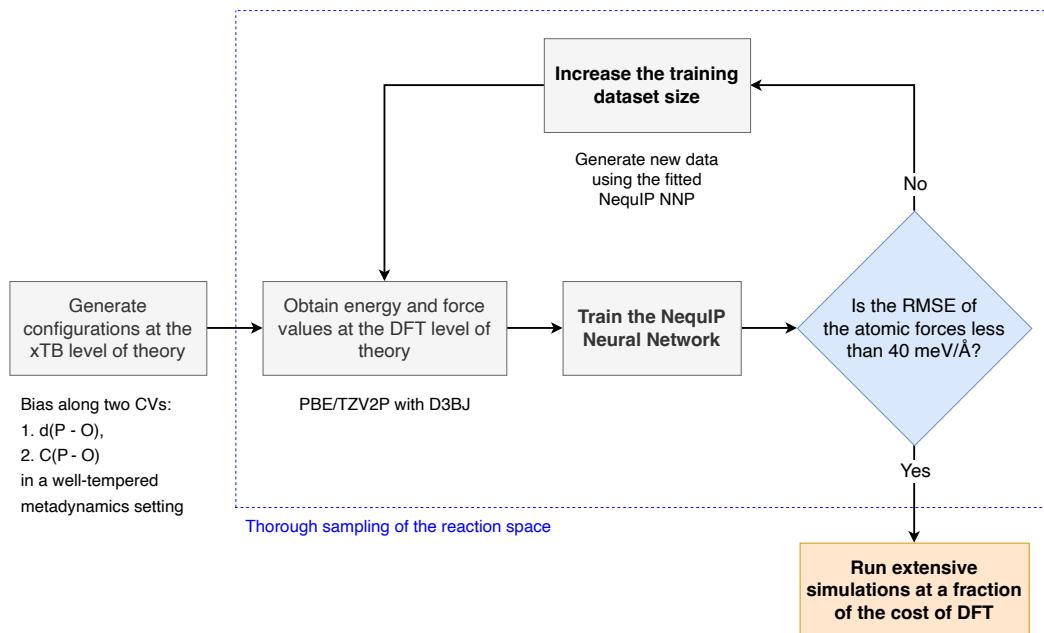


Figure 3.2: Iterative training of the NequIP neural network potential.

each round of training, the model was retrained on a larger dataset. The data obtained from each round will be discussed in the following sections.

Selection of configurations for training and testing

An important part of the iterative training process is the selection of configurations that will be used to train the neural network. To construct a representative and diverse dataset for training the neural network potential, configurations were selected from a metadynamics trajectory using a density-aware sampling strategy. The raw data were extracted from a file generated during the enhanced sampling simulations. Each configuration in this file corresponds to a simulation snapshot, annotated with a time index and two collective variables (CVs): the distance $d(O_{\text{remaining}} - P_{\text{leaving}})$ and the coordination number $\text{CN}(P_{\text{leaving}} - O_{\text{all}})$.

The two CVs were combined into a two-dimensional feature space $\mathbf{X} = (d, \text{CN})$, which served as the basis for sampling. This feature space often exhibits regions of highly non-uniform data density, due to the biased nature of metadynamics sampling. To account for this, a density-aware sampling method was employed to select configurations for training and testing that maintain good coverage across the feature space.

The selection procedure proceeds as follows:

1. A user-defined number of samples is specified.
2. K-means clustering is applied to the feature space to partition it into a number of clusters, k . The number of clusters is determined heuristically as $k = \max(10, \min(\lfloor \frac{N}{50} \rfloor, \lfloor \frac{n_{\text{samples}}}{10} \rfloor))$, where N is the total number of configurations and n_{samples} is the desired number of samples.
3. The number of points sampled from each cluster is proportional to its size, ensuring that denser regions do not dominate the dataset. A minimum of one sample is taken from each non-empty cluster.
4. Within each cluster, a fixed number of configurations is randomly selected using a deterministic random seed to ensure reproducibility.
5. After the training set is selected, the remaining configurations are used to construct the test set, following the same density-aware procedure while ensuring no overlap with the training configurations.

This approach results in training and test datasets that closely mirror the overall distribution of the CVs, while ensuring that underrepresented regions of the feature space are adequately sampled. The final output consists of two lists of snapshot indices corresponding to the selected training and test configurations, along with their respective CV values. These snapshots were then extracted from the trajectory files for use in model training and evaluation. The pseudo-code for the density-aware sampling algorithm is provided in Algorithm A.1.

First round

In the first round of training the NNP, the model was trained on a small dataset consisting of 4,000 configurations. These configurations were obtained from the initial exploration of the configuration space at 300 K using the GFN1-xTB level of theory, as described in Section 3.1.4. The enhanced sampling simulations were biased along $d_{\text{P-O}}$ and $C_{\text{P-O}}$, and no restraints were applied to the system. The training was carried out using the hyperparameters described in Section 3.1.6.

Second round

In the second round of training, the model was trained on a larger dataset consisting of 8,000 configurations. The additional configurations were obtained from a second round of exploration of the configuration space, driven by the NNP obtained after the first round of training.

The NNP-driven simulations were run using the LAMMPS package [51] compiled with PLUMED 2.9.3 [40] and pair_nequip [52] on a single A100 GPU. The simulations were performed for 100 ps in the NVT ensemble at 300 K with the PBC applied in all directions. The temperature was controlled by a Nosé–Hoover thermostat [53, 54] with a time constant of 50 fs. The biasing potential was applied to d_{P-O} and C_{P-O} every 50 fs, using a Gaussian hill height of 2 kcal mol⁻¹ and a width of 0.07 for each CV. The bias factor was set to 30, and the integration time step was 0.5 fs.

Restraints were applied to d_{P-O} and C_{P-O} in order to favour either a dissociative or associative mechanism of the reaction and sample more configurations from the transition state (TS) regions. The training was performed using the same hyperparameters as in the first round.

Third round

In the final round of training, the model was trained on a dataset consisting of 12,000 configurations. These additional configurations were obtained from a third round of exploration of the configuration space, driven by the NNP obtained after the second round of training. The simulations were performed for 500 ps using the same setup as in the second round. The only difference was that the temperature in this round was increased to 320 K and 340 K to explore the configuration space at higher temperatures. The same CVs were biased as in the previous run. No restraints were applied to the system. The training was conducted using the same hyperparameters as in the first round. The final dataset is summarised in Table A.2.

3.2 Production runs at different temperatures

To thoroughly sample the reaction space, production runs were performed at various temperatures (300 K, 320 K, and 340 K) using the neural network potential obtained in the final round of training. To run the simulations with the NNP, the LAMMPS package [51], compiled with PLUMED 2.9.3 [40] and pair_nequip [52], was utilised. First, the systems were equilibrated for 75 ps in NVT ensemble, since it was previously shown that water fully relaxes in this timeframe [55]. Then the simulations were carried out for 4 ns in the NVT ensemble, with temperature regulated by a Nosé–Hoover thermostat [53, 54] with a frequency of 50 fs⁻¹.

The biasing potential was applied every 50 fs to d_{P-O} , C_{P-O} , and C_{O-H} , using a Gaussian hill height of 0.5 kcal mol⁻¹ and a width of 0.07 for each collective variable. Additionally, restraints were imposed on d_{P-O} , C_{O-H} , $C_{P-O_{\text{bridging}}}$, and $C_{P-O_{\text{water}}}$, as mentioned in

Table 3.2. The bias factor was set to 30, and the integration time step was maintained at 0.5 fs. All simulations were conducted on a single A100 GPU.

To obtain the free energy profiles of the reactions, the Gaussian kernels applied during the simulations were summed using the `sum_hills` utility provided in the PLUMED 2.9.3 package [40]. Afterwards, the minimum free energy paths (MFEPs) were extracted using the MEPSA 1.4 software [56].

Having MFEPs in hand and the barrier heights as a consequence, the corresponding rate constants (in s⁻¹) were calculated using the Eyring-Polanyi equation:

$$k = \frac{\kappa k_B T}{h} e^{-\Delta G^\ddagger / RT} \quad (3.4)$$

where $\kappa = 1$. Subsequently, the Arrhenius relationship was derived to obtain the activation energy barrier as the slope of the linear fit from the $\log(k)$ versus $1000/T$ plot.

With the rate constants determined, the corresponding half-lives were then calculated using the following equation:

$$t_{1/2} = \frac{\ln(2)}{k} \quad (3.5)$$

3.3 Validation of the transition states

3.4 Lifetime of the transition states

3.5 Data analysis and visualisation

Chapter 4

Results and discussion

4.1 Final dataset composition

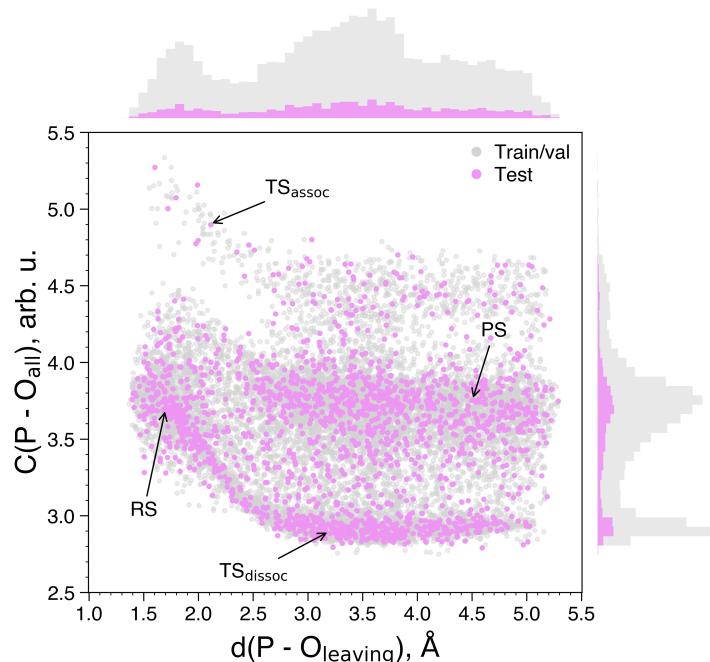


Figure 4.1: Final dataset composition projected on the two CVs space. In total, 12,000 data points are visualised for the training and validation parts, as well as 1,800 points for the test set. RS stands for the reactants state, PS for the product state, and TS for the transition state. 50 bins were used to produce the histograms.

4.2 Accuracy of the neural network potential

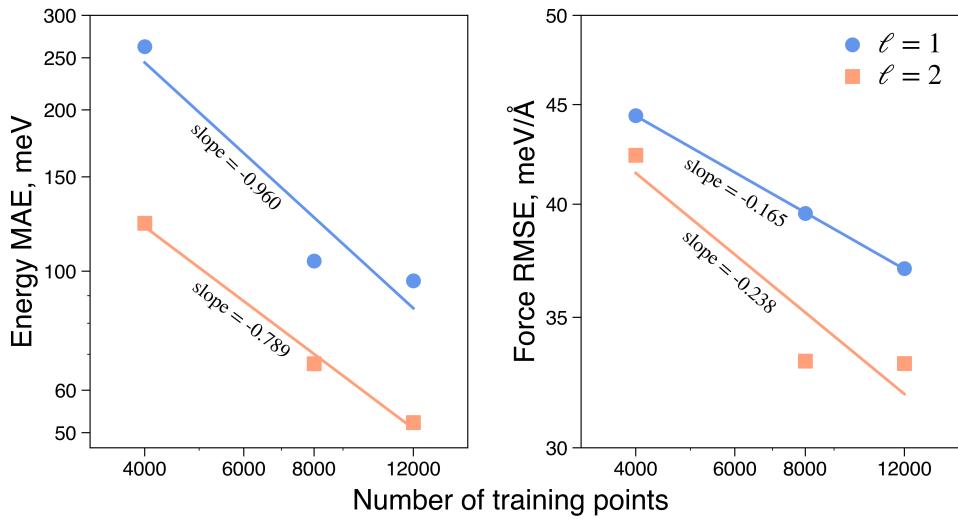


Figure 4.2: Log-log plot of the errors in the energy and forces for the neural network potential with the respect to the training dataset size. In all cases, the errors were calculated on the final test set of 1,200 data points.

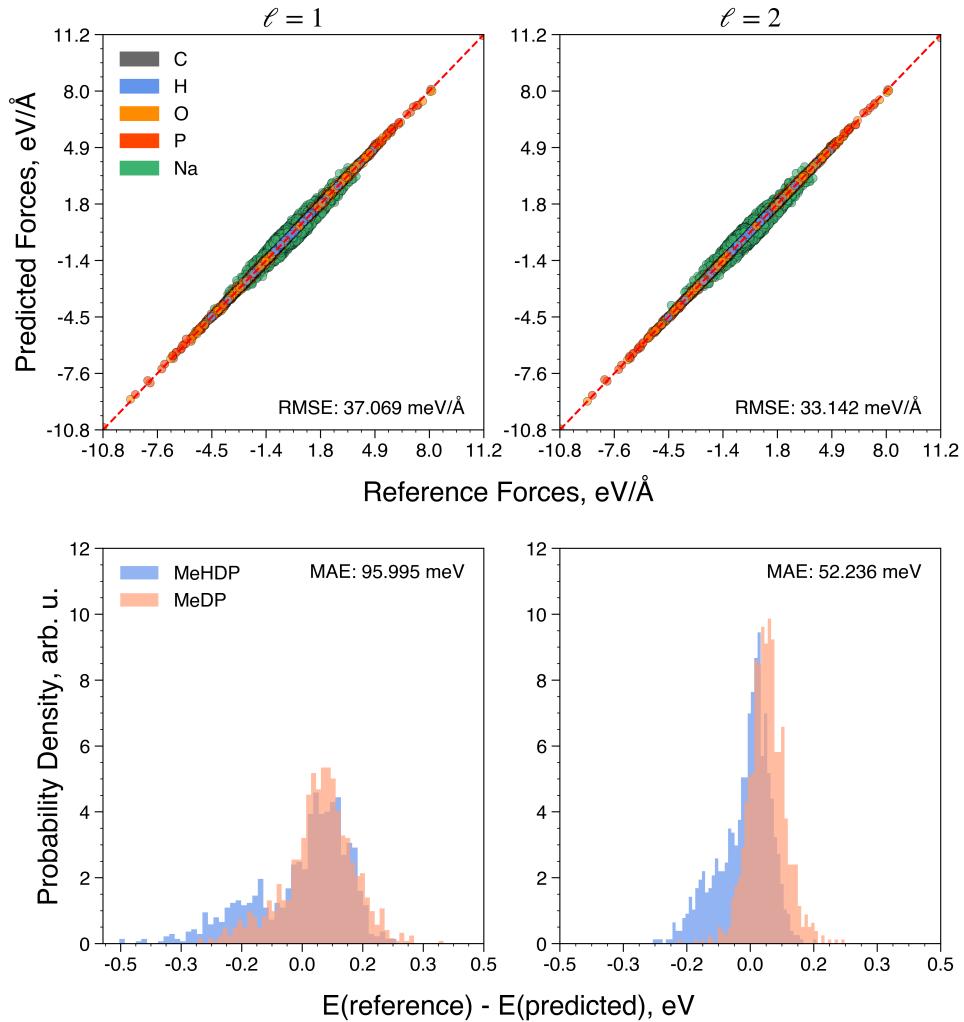


Figure 4.3: Accuracy of the neural network potential trained on 12,000 data points. The left panel shows the errors in the forces and energy for the tensor rank $\ell = 1$ and the right panel shows the errors for $\ell = 2$. The errors are calculated as the difference between the neural network potential and the reference DFT values. For the histograms, the number of bins was set to 50.

4.3 Performance of the potential

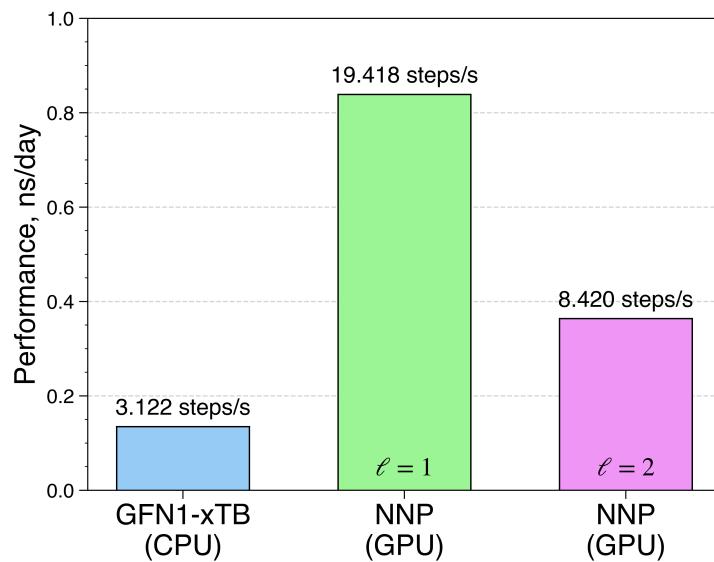


Figure 4.4: Comparison of the performance between the *ab initio* molecular dynamics runs driven by GFN1-xTB and neural network potentials fitted with the different tensor ranks.

4.4 Stability of the production runs**4.5 Convergence of the free energy profiles****4.6 Evolution of the collective variables over time****4.7 Reaction mechanism in case of the methyl diphosphate trianion****4.7.1 Minimum free energy path and free energy surface****4.7.2 Proton transfer mechanism****4.8 Reaction mechanism in case of the methyl diphosphate dianion****4.8.1 Minimum free energy path and free energy surface****4.8.2 Proton transfer mechanism****4.9 Arrhenius relationship**

Chapter 5

Conclusions

Bibliography

- [1] Westheimer, F. H. Why Nature Chose Phosphates. *Science* **235**, 1173–1178 (1987).
- [2] Wolfenden, R. Degrees of Difficulty of Water-Consuming Reactions in the Absence of Enzymes. *Chemical Reviews* **106**, 3379–3396 (2006).
- [3] Müller, W. E., Schröder, H. C. & Wang, X. Inorganic Polyphosphates As Storage for and Generator of Metabolic Energy in the Extracellular Matrix. *Chemical Reviews* **119**, 12337–12374 (2019).
- [4] Nebesnaya, K. S. *et al.* Inorganic polyphosphate regulates functions of thymocytes via activation of P2X purinoreceptors. *Biochimica et Biophysica Acta (BBA) - General Subjects* **1868**, 130523 (2024).
- [5] Kamerlin, S. C. L., Sharma, P. K., Prasad, R. B. & Warshel, A. Why nature really chose phosphate. *Quarterly Reviews of Biophysics* **46**, 1–132 (2013).
- [6] Pavlov, E. *et al.* Inorganic Polyphosphate and Energy Metabolism in Mammalian Cells *. *Journal of Biological Chemistry* **285**, 9420–9428 (2010).
- [7] Dzeja, P. P. & Terzic, A. Phosphotransfer networks and cellular energetics. *Journal of Experimental Biology* **206**, 2039–2047 (2003).
- [8] Boyer, P. D. Energy, Life, and ATP (Nobel Lecture). *Angewandte Chemie International Edition* **37**, 2296–2307 (1998).
- [9] Bonora, M. *et al.* ATP synthesis and storage. *Purinergic Signalling* **8**, 343–357 (2012).
- [10] Walker, J. E. The ATP synthase: The understood, the uncertain and the unknown. *Biochemical Society Transactions* **41**, 1–16 (2013).
- [11] Baev, A. Y. & Abramov, A. Y. Inorganic Polyphosphate and FOF1-ATP Synthase of Mammalian Mitochondria. In Müller, W. E. G., Schröder, H. C., Suess, P. & Wang,

- X. (eds.) *Inorganic Polyphosphates: From Basic Research to Medical Application*, 1–13 (Springer International Publishing, Cham, 2022).
- [12] Baev, A. Y., Angelova, P. R. & Abramov, A. Y. Inorganic polyphosphate is produced and hydrolyzed in F0F1-ATP synthase of mammalian mitochondria. *Biochemical Journal* **477**, 1515–1524 (2020).
- [13] Walker, J. E. ATP Synthesis by Rotary Catalysis (Nobel lecture). *Angewandte Chemie International Edition* **37**, 2308–2319 (1998).
- [14] Watt, I. N., Montgomery, M. G., Runswick, M. J., Leslie, A. G. W. & Walker, J. E. Bioenergetic cost of making an adenosine triphosphate molecule in animal mitochondria. *Proceedings of the National Academy of Sciences* **107**, 16823–16827 (2010).
- [15] Kamerlin, S. C. L., Florián, J. & Warshel, A. Associative Versus Dissociative Mechanisms of Phosphate Monoester Hydrolysis: On the Interpretation of Activation Entropies. *ChemPhysChem* **9**, 1767–1773 (2008).
- [16] Klähn, M., Rosta, E. & Warshel, A. On the Mechanism of Hydrolysis of Phosphate Monoesters Dianions in Solutions and Proteins. *Journal of the American Chemical Society* **128**, 15310–15323 (2006).
- [17] Duarte, F., Åqvist, J., Williams, N. H. & Kamerlin, S. C. L. Resolving Apparent Conflicts between Theoretical and Experimental Models of Phosphate Monoester Hydrolysis. *Journal of the American Chemical Society* **137**, 1081–1093 (2015).
- [18] Hassan, H. A., Rani, S., Fatima, T., Kiani, F. A. & Fischer, S. Effect of protonation on the mechanism of phosphate monoester hydrolysis and comparison with the hydrolysis of nucleoside triphosphate in biomolecular motors. *Biophysical Chemistry* **230**, 27–35 (2017).
- [19] Prasad, B. R., Plotnikov, N. V. & Warshel, A. Addressing Open Questions about Phosphate Hydrolysis Pathways by Careful Free Energy Mapping. *The Journal of Physical Chemistry B* **117**, 153–163 (2013).
- [20] Akola, J. & Jones, R. O. ATP Hydrolysis in Water - A Density Functional Study. *The Journal of Physical Chemistry B* **107**, 11774–11783 (2003).
- [21] Glaves, R., Mathias, G. & Marx, D. Mechanistic Insights into the Hydrolysis of a Nucleoside Triphosphate Model in Neutral and Acidic Solution. *Journal of the American Chemical Society* **134**, 6995–7000 (2012).

- [22] Wang, C., Huang, W. & Liao, J.-L. QM/MM Investigation of ATP Hydrolysis in Aqueous Solution. *The Journal of Physical Chemistry B* **119**, 3720–3726 (2015).
- [23] Harrison, C. B. & Schulten, K. Quantum and Classical Dynamics Simulations of ATP Hydrolysis in Solution. *Journal of Chemical Theory and Computation* **8**, 2328–2335 (2012).
- [24] Ramirez, F., Marecek, J. F. & Szamosi, J. Magnesium and calcium ion effects on hydrolysis rates of adenosine 5'-triphosphate. *The Journal of Organic Chemistry* **45**, 4748–4752 (1980).
- [25] Ramirez, F., , M., James F. & and Szamosi, J. A Comparative Study of Hydrolysis Rates of 2'-Deoxyadenosine and Adenosine 5'-Triphosphates and 5'-Diphosphates. *Phosphorus and Sulfur and the Related Elements* **13**, 249–257 (1982).
- [26] Dirac, P. A. M. & Fowler, R. H. Quantum mechanics of many-electron systems. *Proceedings of the Royal Society of London. Series A, Containing Papers of a Mathematical and Physical Character* **123**, 714–733 (1997).
- [27] Jensen, F. *Introduction to Computational Chemistry* (John Wiley & Sons, 2017).
- [28] Tuckerman, M. E. & Tuckerman, M. E. *Statistical Mechanics: Theory and Molecular Simulation*. Oxford Graduate Texts (Oxford University Press, Oxford, New York, 2023), second edition edn.
- [29] Frenkel, D. & Smit, B. *Understanding Molecular Simulation: From Algorithms to Applications* (Elsevier, San Diego, 2002), second edition edn.
- [30] How the Panama Canal Works. <https://waitbutwhy.com/2014/09/panama-canal-works.html>.
- [31] Jo, S., Kim, T., Iyer, V. G. & Im, W. CHARMM-GUI: A web-based graphical user interface for CHARMM. *Journal of Computational Chemistry* **29**, 1859–1865 (2008).
- [32] Kern, N. R., Lee, J., Choi, Y. K. & Im, W. CHARMM-GUI Multicomponent Assembler for modeling and simulation of complex multicomponent systems. *Nature Communications* **15**, 5459 (2024).
- [33] Kim, S. *et al.* CHARMM-GUI ligand reader and modeler for CHARMM force field generation of small molecules: CHARMM-GUI Ligand Reader and Modeler for CHARMM Force Field Generation of Small Molecules. *Journal of Computational Chemistry* **38**, 1879–1886 (2017).

- [34] Haynes, W. M. *CRC Handbook of Chemistry and Physics* (CRC Press, 2016).
- [35] Bussi, G., Donadio, D. & Parrinello, M. Canonical sampling through velocity rescaling. *The Journal of Chemical Physics* **126**, 014101 (2007).
- [36] Bernetti, M. & Bussi, G. Pressure control using stochastic cell rescaling. *The Journal of Chemical Physics* **153**, 114107 (2020).
- [37] Abraham, M. J. *et al.* GROMACS: High performance molecular simulations through multi-level parallelism from laptops to supercomputers. *SoftwareX* **1–2**, 19–25 (2015).
- [38] Huang, J. *et al.* CHARMM36m: An improved force field for folded and intrinsically disordered proteins. *Nature Methods* **14**, 71–73 (2017).
- [39] Kühne, T. D. *et al.* CP2K: An electronic structure and molecular dynamics software package - Quickstep: Efficient and accurate electronic structure calculations. *The Journal of Chemical Physics* **152**, 194103 (2020).
- [40] Tribello, G. A., Bonomi, M., Branduardi, D., Camilloni, C. & Bussi, G. PLUMED 2: New feathers for an old bird. *Computer Physics Communications* **185**, 604–613 (2014).
- [41] Grimme, S., Bannwarth, C. & Shushkov, P. A Robust and Accurate Tight-Binding Quantum Chemical Method for Structures, Vibrational Frequencies, and Noncovalent Interactions of Large Molecular Systems Parametrized for All spd-Block Elements ($Z = 1\text{--}86$). *Journal of Chemical Theory and Computation* **13**, 1989–2009 (2017).
- [42] Grimme, S., Antony, J., Ehrlich, S. & Krieg, H. A consistent and accurate ab initio parametrization of density functional dispersion correction (DFT-D) for the 94 elements H-Pu. *The Journal of Chemical Physics* **132**, 154104 (2010).
- [43] Barducci, A., Bussi, G. & Parrinello, M. Well-Tempered Metadynamics: A Smoothly Converging and Tunable Free-Energy Method. *Physical Review Letters* **100**, 020603 (2008).
- [44] Perdew, J. P., Burke, K. & Ernzerhof, M. Generalized Gradient Approximation Made Simple. *Physical Review Letters* **77**, 3865–3868 (1996).
- [45] Grimme, S., Ehrlich, S. & Goerigk, L. Effect of the damping function in dispersion corrected density functional theory. *Journal of Computational Chemistry* **32**, 1456–1465 (2011).

- [46] Goedecker, S., Teter, M. & Hutter, J. Separable dual-space Gaussian pseudopotentials. *Physical Review B* **54**, 1703–1710 (1996).
- [47] Hartwigsen, C., Goedecker, S. & Hutter, J. Relativistic separable dual-space Gaussian pseudopotentials from H to Rn. *Physical Review B* **58**, 3641–3662 (1998).
- [48] VandeVondele, J. *et al.* Quickstep: Fast and accurate density functional calculations using a mixed Gaussian and plane waves approach. *Computer Physics Communications* **167**, 103–128 (2005).
- [49] CP2K_Developers. How to Converge the CUTOFF and REL_CUTOFF. <https://manual.cp2k.org/trunk/methods/dft/cutoff.html>.
- [50] Batzner, S. *et al.* E(3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials. *Nature Communications* **13**, 2453 (2022).
- [51] Thompson, A. P. *et al.* LAMMPS - a flexible simulation tool for particle-based materials modeling at the atomic, meso, and continuum scales. *Computer Physics Communications* **271**, 108171 (2022).
- [52] Mir-group/pair_nequip. https://github.com/mir-group/pair_nequip.
- [53] Nosé, S. A unified formulation of the constant temperature molecular dynamics methods. *The Journal of Chemical Physics* **81**, 511–519 (1984).
- [54] Hoover, W. G. Canonical dynamics: Equilibrium phase-space distributions. *Physical Review A* **31**, 1695–1697 (1985).
- [55] Morón, M. C., Prada-Gracia, D. & Falo, F. Macro and nano scale modelling of water–water interactions at ambient and low temperature: Relaxation and residence times. *Physical Chemistry Chemical Physics* **18**, 9377–9387 (2016).
- [56] Marcos-Alcalde, I., Setoain, J., Mendieta-Moreno, J. I., Mendieta, J. & Gómez-Puertas, P. MEPSA: Minimum energy pathway analysis for energy landscapes. *Bioinformatics* **31**, 3853–3855 (2015).

Appendix A

Supplementary information

Table A.1: The plane-wave cutoff convergence test for DFT calculations. The calculation of ΔE involves subtracting the previous energy, e.g. $\Delta E(450 \text{ Ry}) = E(450 \text{ Ry}) - E(400 \text{ Ry})$. When the cutoff ≥ 800 and the rel cutoff ≥ 60 , the error in total energy reduces to ca. 10^{-8} a.u. Only part of the results is shown for the sake of clarity.

Cutoff (Ry)	Rel cutoff (Ry)	Total energy (a.u.)	ΔE (a.u.)
400	60	-2352.6355962810	–
450	60	-2352.6262868887	9.31×10^{-3}
500	60	-2352.6262867349	1.54×10^{-7}
550	60	-2352.6254866602	8.00×10^{-4}
600	60	-2352.6243443853	1.14×10^{-3}
650	60	-2352.6242425582	1.02×10^{-4}
700	60	-2352.6224669798	1.78×10^{-3}
750	60	-2352.6209571227	1.51×10^{-3}
800	60	-2352.6212901605	-3.33×10^{-4}
850	60	-2352.6212901727	-1.22×10^{-8}
900	60	-2352.6212901873	-1.46×10^{-8}
950	60	-2352.6213082173	-1.80×10^{-5}
1000	60	-2352.6208957304	4.12×10^{-4}
10	800	-2354.4562984779	–
20	800	-2352.6775968461	1.78
30	800	-2352.6281701514	4.94×10^{-2}
40	800	-2352.6213637375	6.81×10^{-3}
50	800	-2352.6212892865	7.45×10^{-5}
60	800	-2352.6212901605	-8.74×10^{-7}
70	800	-2352.6212901729	-1.24×10^{-8}
80	800	-2352.6212901739	-1.00×10^{-9}
90	800	-2352.6212901739	0.00
100	800	-2352.6212901739	0.00

Algorithm A.1 Density-aware sampling of configurations**Input:** Feature matrix $\mathbf{X} \in \mathbb{R}^{N \times 2}$ of N configurations, number of samples n_{samples}

1: Determine number of clusters:

$$k \leftarrow \max \left(10, \min \left(\left\lfloor \frac{N}{50} \right\rfloor, \left\lfloor \frac{n_{\text{samples}}}{10} \right\rfloor \right) \right)$$

2: Apply K-means clustering to \mathbf{X} with k clusters3: Initialize empty list for sampled configuration indices $S \leftarrow []$ 4: **for** each cluster C_i , $i = 1$ to k **do**

5: $n_i \leftarrow \max \left(1, \left\lfloor \frac{|C_i|}{N} \cdot n_{\text{samples}} \right\rfloor \right)$

6: Select n_i random configurations from C_i with fixed random seed7: Append selected indices to S 8: **end for**9: **return** S **Output:** List of selected configuration indices S

Table A.2: Composition of the full dataset used for training and testing. Well-tempered meta-dynamics settings used to run the simulations: ¹GPN1-xTB for energies and forces, gaussian height = 2 kcal/mol, spawning frequency = 25 fs⁻¹, bias factor = 30 and ²NNP for energies and forces, gaussian height = 2 kcal/mol, spawning frequency = 50 fs⁻¹, bias factor = 30.

System	Temperature (K)	Simulation length (ps)	Train/Val	Test
MeDP ¹	300	50 ps	2000	150
MeDP ²	300	100 ps	2000	150
MeDP ²	320	500 ps	1000	300
MeDP ²	340	500 ps	1000	300
MeHDP ¹	300	50 ps	2000	150
MeHDP ²	300	100 ps	2000	150
MeHDP ²	320	500 ps	1000	300
MeHDP ²	340	500 ps	1000	300
Total			12000	1800

Quantum Chemistry and Physical Chemistry

Celestijnenlaan 200F bus 2404

3001 LEUVEN, BELGIË

tel. + 32 16 37 21 98

jeremy.harvey@kuleuven.be

www.kuleuven.be

