

# Computer Modelling the Root Cause of Cystic Fibrosis

by

Miro Alexander Astore

*A thesis submitted in fulfilment of the  
requirements for the degree of*

Doctor of Philosophy

School of Physics  
Faculty of Science  
The University of Sydney

2022

Declaration of Original contribution

of the dissertation submitted by

Miro Alexander Astore

This is to certify that to the best of my knowledge, the content of this thesis is my own work. This thesis has not been submitted for any degree or other purposes.

I certify that the intellectual content of this thesis is the product of my own work and that all the assistance received in preparing this thesis and sources have been acknowledged.

---

*Miro Alexander Astore*, Author

---

Date

## **Abstract**

placeholder text

*In loving memory of Madeline Jennifer Dell*

“Fear cuts deeper than swords.”

Arya Stark

## Acknowledgments

Daniel Golestan, a wise man, once told me that to be given the opportunity to create this thesis was a gift. It was. It was a gift given to me by every friend, colleague, teacher, mentor and family member I've spent any time with. The list that follows of those to thank is not complete. If it was you'd be reading about a conversation I had with a middle aged woman in a hostel north of San Francisco, but that has little to do with Cystic Fibrosis.

To My parents raised me with not only academic rigor in mind but also a respect for aesthetics which has served me strangely well. I've never had a talent for the creative side of things compared to quantitative disciplines. But were it not for their demand for respect for the arts I'd have remained illiterate.

To Jeffry for his tutelage and patience, even across the pacific ocean. To have been your first mentee is an honor. You will go far.

To Poker, I am a better human being in every conceivable way for having known you. Your wisdom, intelligence and kindness are boundless. You have taught me an inordinate number of things. And yes, I do mean inordinate.

Nono and Nona I don't think you'll ever read this. I'm sad that you won't understand what I've done but I think you'd be proud if you did. Living in Condell park did more for me than you could know. Far from war torn Beirut or dirt poor Orria I'm sitting in a well lit office writing this with a full stomach and few worries. Sometimes this luck makes my head spin.

Thank you to Shafagh Waters for her vision, her drive and all her advice. You brought me a truly fascinating PhD project and I benefited greatly from your mentorship. Bridging the gap between cell biology and molecular physics is something that will happen more in the future and I'm lucky to have met such a driven lab to teach me to do so.

Serdar, a brilliant mind and a patient boss. Thank you for giving me the best possible experience at grad school I could have asked for. Your willingness to let me pursue self directed projects with a guided hand is a privilege during a PhD and I'm all the better for having gotten it from one of the best. I'm excited to carry some of your physical insight into biological systems to future research projects.

Maddy, I miss you every day. You couldn't have imagined what it was like to do this after losing you. I carry much of you with me and I wish I had more. I miss your intelligence, your warmth and your love.

You're all in my Loop and I hope I'm in yours in some way.

## List of Publications

MA - Miro Alexander Astore

SK - Serdar Kuyucak

1. placeholdertext

## Publication Authorship Attribution

In addition to the statements above, in cases where I am not the corresponding author of a published item, permission to include the published material has been granted by the corresponding author.

---

*Miro Alexander Astore*, Student

Date

As the supervisor for the candidature upon which this thesis is based, I can confirm that the authorship attribution statements above are correct.

---

*Serdar Kuyucak*, Supervisor

Date

# **Contents**

# List of Abbreviations

<i>AMBER</i>	Assisted Model Building with Energy Refinement
<i>BAR</i>	Bennett-Acceptance-Ratio
<i>CF</i>	Cystic Fibrosis
<i>CFTR</i>	Cystic Fibrosis Transmembrane Conductance Regulator
<i>CHARMM</i>	Chemistry at Harvard Macromolecular Mechanics
<i>COM</i>	Centre of Mass
<i>CV</i>	Collective Variable
<i>FEP</i>	Free-Energy Perturbation
<i>gA</i>	Gramicidin A Ion Channel
<i>Glt<sub>Ph</sub></i>	Glutamate Transporter - <i>Pyrococcus horikoshii</i>
<i>GROMACS</i>	GROningen MAchine for Chemical Simulations - MD program
<i>GROMOS</i>	GROningen MOlecular Simulation - MD program
<i>LJ</i>	Lenard-Jones Potential
<i>MBAR</i>	Multistate Bennett-Acceptance-Ratio
<i>MD</i>	Molecular Dynamics
<i>MetaD</i>	Meta Dynamics
<i>NAMD</i>	Nanoscale Molecular Dynamics - MD Program
<i>NBD</i>	Nucleotide Binding Domain
<i>NPT</i>	Constant number of Particles, Pressure and Temperature
<i>NVT</i>	Constant number of Particles, Volume and Temperature
<i>OpenMM</i>	Open Molecular Mechanics - MD Program
<i>OPLS</i>	Optimised Potentials for Liquid Simulations
<i>PBC</i>	Periodic Boundary Condition
<i>PCA</i>	Principal Component Analysis
<i>PDB</i>	Protein Data Bank
<i>PMF</i>	Potential of Mean Force
<i>PME</i>	Particle Mesh Ewald - Long-range Electrostatics Method
<i>POPC</i>	1-palmitoyl-2-oleyl-sn-glycero-3-phosphocholine
<i>POPE</i>	1-palmitoyl-2-oleyl-sn-glycero-3-phosphoethanolamine
<i>RMSD</i>	Root-Mean-Square Deviation
<i>TI</i>	Thermodynamic Integration
<i>TICA</i>	Time-lagged Independent Component Analysis
<i>US</i>	Umbrella Sampling
<i>VMD</i>	Visual Molecular Dynamics - MD Visualisation Program
<i>WHAM</i>	Weighted Histogram Analysis Method



# List of Figures

# List of Tables





# Chapter 1

## Introduction

*Whatever complexity means, most people agree that biological systems have it. -Frauenfelder and Wolynes??*

### 1.1 Physics in a test tube

Why can't I write down an equation will tell me how long I will live? Or how many hairs I will grow?

This might seem like an inane question but if you asked a physicist for the formula for how long it takes a radioactive material to decay or how long it will take an object to fall into a black hole they will be able to answer easily.

What makes the first set of questions so much more difficult to answer?

I posit that it is the diversity of components that makes biological questions so difficult to ask and answer. Biology distinguishes itself amongst scientific disciplines requiring the study of systems that are both complex and heterogeneous. In the study of more simple physical systems a simple analogy such as a mass on a spring or a gas of hard spheres can be extremely successful in explaining macroscopic phenomena. For biological systems there appears to be too much complexity for such analogies to have the same level of success. They may struggle to answer questions such as "If this gene mutates how will that affect lung function?" "If this drug were given at a higher dosage what would its effect be?" "What if we change this chemical moiety?" At the moment, a trained chemist needs to go and answer these questions pipette in hand, the physicist with their notebook is hopeless.

It seems like a silly question but it seems important to ask why we can't just use a device similar to a harmonic oscillator or a perfect black body to speculate at useful

answers for these quantitative questions. The answer is just as silly. If you look with your naked eye at your arm, you will notice hair, pores, dry skin, dead skin, perhaps even tendons and muscles under the skin. If you take a microscope you will notice the 3 layers to your skin with different functions and composition. If you were to take a single cell from any of those layers and stain it to distinguish features in an electron microscope you would notice all sorts of complex structures and the size and number of these structures would vary depending on where you took the cell from in the body. Within and between each those structures is a salty, wet dance of molecules large and small. This heterogeneity on length scales hints at the reasons behind biology's physical complexity. Plasma physics is often characterised by the density of the plasma studied. This parameter may span 28 orders of magnitude from a dense stellar core to the sparse intergalactic nebulae. The same mathematical tools can be used to map any plasma in these energy scales. Would that we were so lucky in biology. We struggle to apply same physical models to deal with phenomena across a single order of magnitude.

Thus, in order to move towards more predictive theories of biology it is necessary to consider much more of the fundamental physical processes occurring within biological systems than simply searching for statistical trends. One form of this from fundamentals approach is the simulation of every atom in a biological system. Although computationally expensive, this approach appears necessary due to the heterogeneous nature of biological systems.

One of the things we're trying to do with molecular dynamics is fill in the gap left by the sequence- $\rightarrow$ function paradigm which is internalised in current understandings of molecular biology. We usually talk about how the sequence of the gene defines its function because it gives the protein its structure but really there is a considerably larger amount of regulatory pressure exerted by the environment. This is what is missing from the sequence alone paradigm.

## 1.2 What is Physics?

Personally I have always given answers along the lines of "the study of the movement of energy within a system" or when I was in high school "The study of how things move". Although adequate for a layman these might obscure the fundamental structure within physics that make it such a powerful tool. It is the conception of some causal unit in a system and the ability to scale up the behaviour of that unit to make predictions about measurable phenomena.

This might take a few different forms at different scales, it's what makes physics feel like the most "fundamental" of the sciences.

Examples include:

Newton's laws of gravitation to explain the organisation of the solar system.

Einstein's theories employing Riemannian geometry to track the motions of galaxies and black holes.

The conception of atoms as hard spheres used to derive the macroscopic behaviour of

gasses.

The schrodinger wave function to find the structure of atoms, which can then be integrated further up to find their macroscopic organisations. More on this later.

Biological systems exhibit such a problem for the physicist because unlike the above problems it is extremely hard to pick out a fundamental unit to even begin our upwards journey. An evolutionary biologist might say to choose the "gene" but this is actually far too high in our spatial heirarchy already. Really a gene is only meaningful to the dance of life if it has partners to dance with. Genes of hard spheres ?

A coil of DNA in water doesn't really do much in solution except decay without machinery that can preserve, read, translate and replicate it. The gene is an emergent property, we have to go deeper.

So, what creates the gene?

A slew of biological machinery that mostly take the form of proteins. These proetins are a special case of chemistry, with many observable functions. Their sequence is coded by the DNA in something reminiscent of a strange loop [Hoffstadter2008].

This self referential loop is one of the reasons biology is so difficult. Since we know that this strange loop is kicked off by atomic interactions we will start there. As we are taking a physical, pragmatic approach here it would make sense to begin with the protein, after all, they stave off the march of entropy constantly trying to eat up all of your cells. It also just so happens that they are much easier to understand computationally since their motions are faster and more flexible.

The first level sub cellular organisation is perhaps the most intimidating first step for me personally after spending 4 years simulating a single protein. Glimpsing the complexity within a single one of these molecules has been one of the most existential experiences of my life but the knowledge that there are astronomical numbers of these things inside me all of the time

It is hoped that illustrating the monumental task in both intellectual effort and resources of incrementally increasing the understanding of a single protein amongst tens of thousands will give the reader and understanding of how we might continue our quest to understand the molecular dance that plays within all of us.

This makes sense if we think about it Somewhere on the scale between a single protein and a single cell this is what we consider "life". We have single unicellular organisms but we don't have uniproteomic organisms. So the fundamental length scale of life is somewhere between  $10^{-10}m$  and  $10^{-3}m$ . This is the first loop in our strange loop.

After this things start to run away from me with my handful of GPUs and limited patience. So in this thesis we will only discuss single proteins.

## 1.3 Ion Channels: Natures laboratories to Teach Us Biophysics

The physiological importance of ion channels became clear after the experiments of Hodgkin and Huxley. These mathematicians took nerves from fished giant squid and measured the current running through the nerve in response to electrical stimulation. What they found was intriguing. Current would only flow when the input signal was of a sufficient voltage. The measurements and modelling they carried out gave an exciting set of results. They found that the cell had to maintain a constant electrochemical gradient, they discovered that the presence of voltage gated ion channels and cation selective ion channels[1]. Each of these features, motivated by mathematical modelling have been found to be critical to the functioning of the cell and fundamental to the foundation of molecular biophysics. The following set coupled ordinary differential equations were discovered by testing functions which fit the measurements taken from the squid axon.

$$\begin{aligned} I &= C_m \frac{dV}{dt} + \bar{g}_K n^4 (V - V_K) + \bar{g}_{Na} m^3 h (V - V_{Na}) + \bar{g}_l (V - V_l), \\ \frac{dn}{dt} &= \alpha_n(V)(1 - n) - \beta_n(V)n, \\ \frac{dm}{dt} &= \alpha_m(V)(1 - m) - \beta_m(V)m, \\ \frac{dh}{dt} &= \alpha_h(V)(1 - h) - \beta_h(V)h \end{aligned} \tag{1.1}$$

The  $\alpha$  and  $\beta$  parameters are the proportion of the sodium and potassium channel populations which are activated, respectively. This example shows how basic theoretical tools can be used to predict and discover physical phenomena in biological systems. The Hodgkin Huxley model proved the existence of a cell's resting potential, the possibility of voltage gated ion channels, and channels whose pores are selective for certain ions. Even today the molecular mechanisms behind some of these discoveries are debated. In this thesis we aim to do the same by building up from fundamental quantum mechanics in order to understand the motion of single proteins so we might speculate as to the function of the whole organism.

Similar to the above story, ion channels have always motivated the early pioneers of molecular biophysics. This is due to their ubiquity and importance in biological systems and the ease of measuring their activity with biochemical assays. One just needs an oscilloscope to measure their current. As cell biology has advanced it has become clear that the resting potential of a cell is critical to its function, regulating many chemical reactions inside it.

These factors have allowed biophysicists sufficient data to build sufficiently accurate models of biomolecular systems which generalise to other systems. Leading to a thriving field, analysing systems as diverse as protocells to gold nano particles CITATIONS NEEDED.

The discovery of voltage gated channels and a resting potential .

## 1.4 Studying Cystic Fibrosis to Learn Biophysics

The sad truth of this debilitating disease is that those afflicted are extremely unlucky. A single, small change to the genome and their lungs fill with sticky mucus and become infected with bacteria, each breath cumbersome. Personally, I've not met somebody who has this disease. I have consistently wondered what perspective I'm missing by not suffering myself from such a condition or even knowing somebody with it. I'm not been trained in the ethics of studying medicine.

In this way, my motivations for studying this protein aren't solely focussed on treating disease. There is a perspective on protein evolution which states that the primary sequence of a particular gene contributes to the overall fitness of an organisms by a formula. []

It just so happens that the CFTR gene sits at the precipice of a daunting cliff in sequence space. So by taking small steps in sequence space and plunging down this cliff we can try to understand how we might push the ball back up the cliff and retain functionality.

Moreover, by learning the nuts and bolts of what goes wrong with CFTR we can start to think about where some of these cliffs might be in other places in the proteome, to gain function and avoid disease and debilitation.

The reality of disease pathogenesis being caused by so many different mutations means that there has been decades of investigation into the function of every domain in the protein.

Due to the array of disease causing mutations which occur accross the cystic fibrosis protein, there is a large body of literature on its unique function. This allows us a glance into its function and an opportunity to simultaneously perform basic biophysical research while directly assisting in furthering patient outcomes.

## 1.5 Well. We're in the future

Throughout science, the integration of experimental data with theoretical models leads to new and exciting research, this is particularly true in biology with its important applications in medicine, agriculture and manufacturing. Wet lab biologists take advantage of experimental techniques which allow them to understand the dynamics and structure of living things from the top down. The finer the experimental instrument, the finer the detail they may resolve. Conversely, computational and theoretical biologists take a bottom up approach, we aim to take the granular details of a system, and integrate them upwards to model the macroscopic behaviour of that system. With more powerful computers and more detailed models we can make predictions about the behaviour of more complex systems. What is so exciting about the current era of biological research is that the domains of these two approaches are beginning to overlap, where they can synergize and drive further breakthroughs. As we discover more systems where this overlap can be found we will develop more sophisticated treatments for diseases and problems found around the world.

The reason this has happened before in physics is two fold. Physical systems are much more homogeneous. So it's much easier to integrate upwards in length scale. Once you understand the pairwise interaction between two components it's simply a question of having the theoretical and computational capacity to model the bulk behaviour of that system.

The difference with biological systems is that they have so many different components that finding an analytic or even computationally tractable solution is usually impossible. However, as we collect more data and build more powerful computers we can approach more complete models. These in turn inform more powerful theoretical models these help direct the material efforts of experimental expertise .

AlphaFold is a good example. This new breakthrough builds on decades of inquiry from the structural biology community and advancements in AI to give high resolution protein structures. Now this result can be used to fill in the gaps of structural biology. Crucially, AlphaFold knows what it doesn't know. So we can tell where to direct the efforts of structural biology. Together these advances will fill more gaps in our knowledge of protein physics.

# Chapter 2

## From Protons to Proteins: Methods to simulate the inside of a cell.

The purpose of this chapter is to train those who have studied physics in some of the details they will need to understand the models we use to simulate molecular systems (and the many technical problems they will encounter). An excellent overview which I would recommend as first reading for any new student can be found in an article by Braun et al. [2]. We will go flesh out the physics in some more detail here but this article provides a broader overview of different techniques and resources for where one might be able to find more physical details. We assume natural units for all equations.

### 2.1 Quantum Mechanics is Not Tractable at the Scale of Biology.

Living things are made of atoms and atoms themselves are composed of many particles. The motions of atoms and their constituent particles are governed by quantum mechanics. Unfortunately, performing simulations for the number of atoms involved in proteins and other cellular components at quantum mechanical accuracy is impossible. Hence, we will show how to take the fundamental formulation of atomic interactions in the Schrödinger wave equation and apply approximations in order to produce a model which is capable of simulating macromolecular systems at biologically relevant timescales.

We will gradually integrate upwards, beginning with the interactions in a single atom we will work our way up to a complex macromolecular system with lipids, water, salts and of course, proteins. Ultimately this section rationalises the treatment of atoms as point charges in classical molecular dynamics simulations.

#### 2.1.1 A full quantum mechanical treatment

Since we are dealing with atoms which are governed by quantum mechanics we must begin our journey upwards with the time dependent form of the Schrödinger wave

equation.

$$i\hbar \frac{\partial}{\partial t} \Psi(\mathbf{x}, t) = \left[ -\frac{\hbar^2}{2m} \nabla^2 + V(\mathbf{x}, t) \right] \Psi(\mathbf{x}, t) \quad (2.1)$$

In quantum systems we treat all particles as waves hence the use of the wave function  $\Psi(\mathbf{x}, t)$ . The complex amplitude of the wave function  $|\Psi(\mathbf{x}, t)|^2$  tells us the likelihood of detecting the particle at time  $t$  and at place  $\mathbf{x}$ . The term in the brackets correspond to  $-\frac{\hbar^2}{2m} \nabla^2$  the kinetic energy of the particle with mass  $m$  while  $V(\mathbf{x}, t)$  is the potential energy of the system. Given that the left hand term  $i\hbar \frac{\partial}{\partial t} \Psi(\mathbf{x}, t)$  contains a gradient with respect to time, it governs how the wave function will evolve in time.

When the external potential  $V$  has no explicit dependence on time, this equation reduces to the familiar time independent form.

$$E\Psi(\mathbf{x}, t) = \left[ -\frac{\hbar^2}{2m} \nabla^2 + V(\mathbf{x}) \right] \Psi(\mathbf{x}, t) = H\Psi(\mathbf{x}, t) \quad (2.2)$$

Note that the wave function  $\Psi(\mathbf{x}, t)$  is still allowed to evolve in time.

In atomic systems there are two types of particles, nuclei which we will denote with the subscript  $i$  and electrons denoted by  $e$ . In order to treat these elements separately we decompose the Hamiltonian of the system into a few components.

$$H = \underbrace{T_n + V_{n-n}}_{H_n} + \underbrace{T_e + V_{e-e} + V_{n-e}}_{H_e} \quad (2.3)$$

Where  $T_n$  and  $T_e$  denote the kinetic energy of the nuclei and electrons respectively. While  $V_{n-n}$ ,  $V_{n-e}$ ,  $V_{e-e}$  denote the potential energy for interactions between nuclei, between electrons and nuclei and between electrons respectively.

Since the potential terms all describe charged species, they follow Coulomb's law and have the form.

$$V_{n-n} = \sum_{i>j} \frac{q_e^2 z_i z_j}{|\mathbf{R}_i - \mathbf{R}_j|}, \quad V_{n-e} = - \sum_{i,l} \frac{q_e^2 z_i}{|\mathbf{r}_l - \mathbf{R}_i|}, \quad V_{e-e} = \sum_{l>k} \frac{q_e^2}{|\mathbf{r}_l - \mathbf{r}_k|} \quad (2.4)$$

Here the  $z_i$  represent the atomic number (and thus the charge) of the  $i$ th nucleus and  $q_e$  is the unit charge of the electron. The reason for the separate coordinates  $R_i$  and  $r_l$  is to separate out the treatment of nuclei and electrons which will be important once we apply the Born-Oppenheimer approximation.

Meanwhile, the kinetic energy terms are of the form

$$T_n = - \sum_i \frac{\hbar^2}{2M_i} \nabla_i^2, \quad T_e = - \sum_l \frac{\hbar^2}{2m_e} \nabla_l^2 \quad (2.5)$$

$M_i$  represents the mass of the  $i$ th nucleon and  $m_e$  represents the mass of an electron. The operator  $\nabla^2 = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2}$ . The separate subscripts  $i$  and  $l$  are due to the different coordinates which we use to denote the positions of the nuclei and the electrons. The reason for this will become clear when we derive the Born-Oppenheimer approximation to separate the wave functions and treat them separately.

### 2.1.2 The Born-Oppenheimer approximation.

In order to reach the Born-Oppenheimer approximation we start with the observation that electrons have a mass 3-4 orders of magnitude smaller than the nuclei. This motivates two simplifications. The ‘clamped nuclei assumption’ where we solve the Schrodinger equation whilst nuclei are fixed in space and do not move. And a related assumption known as the “adiabatic assumption” which postulates that the electrons will respond instantaneously to any changes in the positions of the nuclei. Combining these physical approximations we derive the “Born-Oppenheimer approximation” for the Schrodinger equation which can be used to simplify calculations involving several atoms at once.

We begin the derivation by examining the time-independent form of the electronic Schrodinger wave equation where the nuclei are fixed at positions  $R_i$ .

$$H_e(\mathbf{r}_l, \mathbf{R}_i)\psi_e(\mathbf{r}_l, \mathbf{R}_i) = E_e(\mathbf{R}_i)\psi_e(\mathbf{r}_l, \mathbf{R}_i) \quad (2.6)$$

Fixing the nuclei in this way gives the “clamped nuclei” approximation [3]. To solve the wave function for the whole system  $\Psi_{tot}$  we use an *ansatz* which decomposes the wave function with an electronic basis into two components:  $(\psi_e)_k$  and  $(\psi_n)_k$  which are the  $k$ th eigenfunction solutions to  $H_e$  and  $H_n$  respectively.

$$\Psi_{tot}(\mathbf{r}_l, \mathbf{R}_i, t) = \sum_{k=0}^{\infty} \psi_e(\mathbf{r}_l, \mathbf{R}_i)_k \psi_n(\mathbf{R}_i)_k \quad (2.7)$$

Note that there is an implied direct product between the wave functions  $\psi_e(\mathbf{r}_l, \mathbf{R}_i)$  and  $\psi_n(\mathbf{R}_i)$ . When we substitute this expression into the full Schrodinger equation 2.1 we find the following expression for the  $k$ th nuclear eigenfunction [4]

$$i\hbar \frac{\partial}{\partial t} \psi_n(\mathbf{R}_i)_k = \left[ - \sum_i \frac{\hbar^2}{2M_i} \nabla_i^2 + E_e(\mathbf{R}_i)_k \right] \psi_n(\mathbf{R}_i)_k + \sum_j C_{kj} \psi_n(\mathbf{R}_i)_j \quad (2.8)$$

Where we have coupled the electronic wave functions to each other with the operator

$$C_{kj} = \int (\psi_e)_k^* \left[ \sum_i \frac{\hbar^2}{2M_i} \nabla_i^2 \right] (\psi_e)_j d\mathbf{r} + \frac{1}{M_i} \sum_i \left[ \int (\psi_e)_k^* [-\hbar i \nabla_i] (\psi_e)_j d\mathbf{r} \right] [-\hbar i \nabla_i] \quad (2.9)$$

Using the “adiabatic assumption” [4] the off-diagonal terms of  $C_{kj}$  can be set to 0. This completely decouples the wavefunction into two components

$$\Psi_{tot}(\mathbf{r}_l, \mathbf{R}_i, t) = \psi_e(\mathbf{r}_l, \mathbf{R}_i)_k \psi_n(\mathbf{R}_i, t)_k \quad (2.10)$$

Since all cross terms from the direct product can be ignored. With the further assumption that the diagonal terms  $C_{kk}$  can also be ignored because they are 4 orders of magnitude smaller than the other terms in 2.8 [3].

We now write the Born-Oppenheimer approximated wave equation for an atomic system.

$$i\hbar \frac{\partial}{\partial t} \psi_n(\mathbf{R}_i)_k = \left[ -\sum_i \frac{\hbar^2}{2M_i} \nabla_i^2 + E_e(\mathbf{R}_i)_k \right] \psi_n(\mathbf{R}_i)_k \quad (2.11)$$

By rearranging this equation and taking derivative we can see how to use Newton’s equations of motion to calculate the forces on the nuclei from the surrounding electric potential

$$M_i \ddot{\mathbf{R}}_i(t) = -\nabla_i E_e(\mathbf{R}_i) \quad (2.12)$$

By choosing an appropriate time-step one can simply iteratively solve this equation of motion to understand the dynamics of an atomic system. The nuclei will move according to their relative positions to each other and the electron clouds will rearrange in response to that motion. There is no need to calculate the interactions between the nuclei and the electrons. This is sufficient accuracy to simulate many physical phenomena with the notable exception of energetic, fast interactions between nuclei and electrons such as spectroscopic phenomenon [3].

## 2.2 Classical MD, Molecular Motions Without Quantum Mechanics

The Born-Oppenheimer approximation gives rise to Hartree-Fock methods and density functional theory (DFT). These more sophisticated physical methods allow us to simulate the organisation of electron clouds around small molecules, finding broad applications in chemistry and materials science [5]. We can derive the energy profile of certain degrees of freedom within the molecule such as the energetics of stretching out a bond or twisting a dihedral angle. These can be useful when designing novel materials.

However, even with these approximations simulating a large number of atoms is still not computationally tractable. State of the art DFT methods can only simulate up to a few 10s of thousands of atoms [6] and scales as  $O(N^3)$  [7]. This is not sufficient to simulate proteins and their surrounding solvation environment. So, we must use another round of approximations to reach the spatial and time scales necessary to

simulate biological molecules. We do this by creating a set of mathematical functions the calculations further. Here we use a set of virtual springs and other simple models for the energetic interactions between atoms. This creates what's known as an effective potential. So named because it effectively approximates the behaviour of the full quantum mechanical system.

This formulation gives us classical molecular dynamics sometimes referred to as molecular mechanics. The aim of the classical forcefields discussed here is to use *ab initio* MD as a target to approximate.

The CHARMM effective potential employed in this work is similar to those found in all common all-atom molecular dynamics forcefields. The same functional forms are used in other forcefields such as AMBER, GROMOS and OPLS but with different parameters and design philosophies. [CITATION NEEDED]

We split up the molecular potential into several components dealing with the energies from covalent bonds, including bond stretching, twisting and pinching. As well as energies associated with the forces that atoms exert on each other when they are not bonded together. Namely Coulomb forces due to electric charges on the atom and attractive Van Der Walls interactions and repulsion due to Pauli Exclusion the latter two forces are combined into one term we will analyse in detail  $U_{LJ}$ .

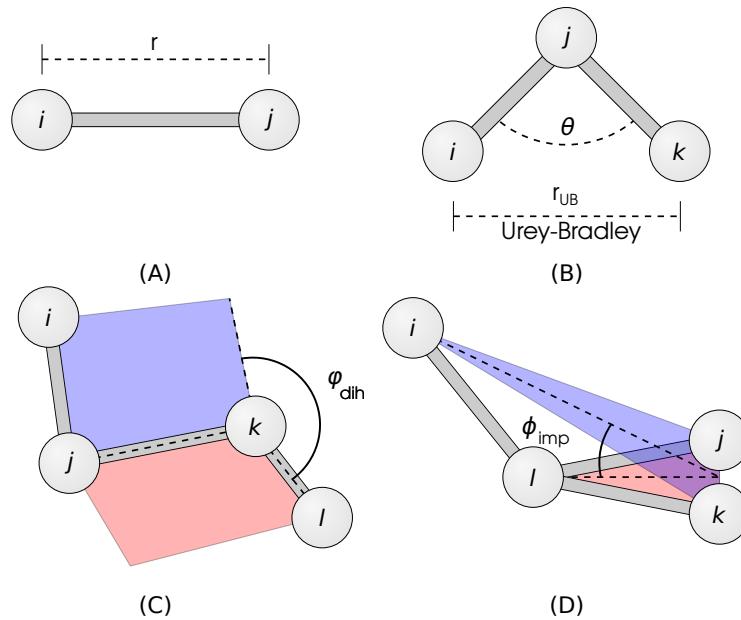
$$U_{MM} = \underbrace{U_{LJ} + U_{Coulomb}}_{U_{non-bonded}} + \underbrace{U_{bonds} + U_{angles} + U_{dihedrals} + U_{impropers}}_{U_{bonded}} \quad (2.13)$$

Interestingly, the bonded terms may all reasonably be approximated by harmonic springs.

$$\begin{aligned} U_{bonded} = & \sum_{bonds} k_b(b - b_0)^2 + \sum_{angles} k_\theta(\theta - \theta_0)^2 + \sum_{Urey-Bradley} k_u(r_{UB} - r_{UB_0})^2 \\ & + \sum_{dihedrals} k_\varphi(1 + \cos(n\varphi - \delta)) + \sum_{improper-dihedrals} k_\phi(\phi - \phi_0)^2 \end{aligned} \quad (2.14)$$

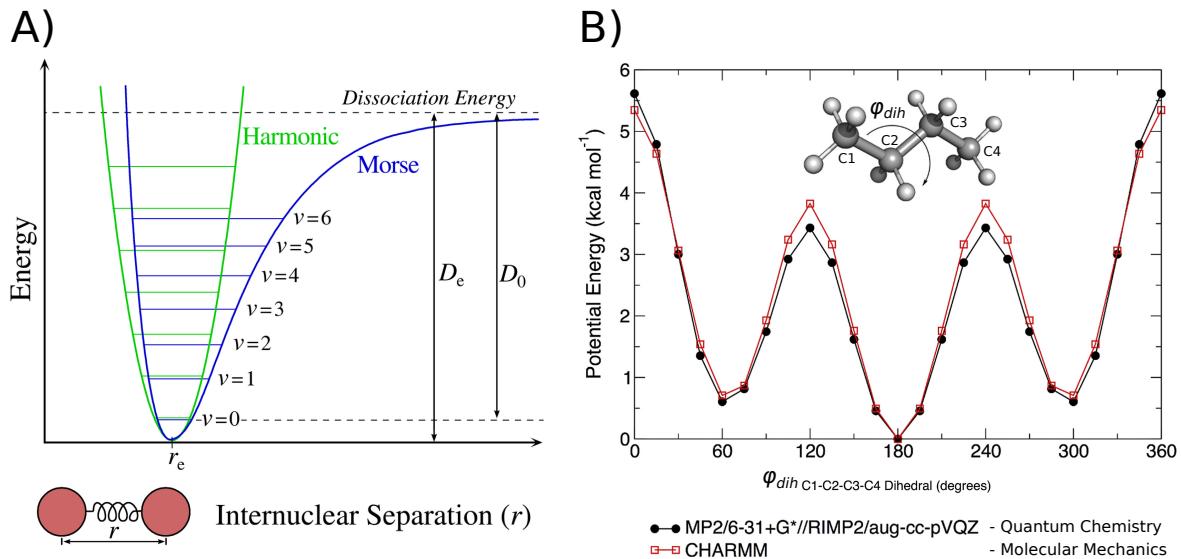
Here, the  $k_i$  terms correspond to the strength of the harmonic restraint for that parameter. The 0 subscript denotes the equilibrium position for that parameter. Even though this formulation is quite simple, it has empirically been shown to be a reasonable approximation for the potential energy functions of quantum mechanics in covalently bonded bonded species. Examples can be seen in figure 2.2.

In classical forcefields the non-bonded interactions are expressed using the Couloumb's law because the partial charges assigned to each atom and the Lennard-Jones potential to approximate the interactions arising from both Pauli exclusion and Van Der Walls Interactions.



**Figure 2.1: The Bonded Interactions Calculated In Classical Forcefields.**

(A) The energy of Bond Stretching is approximated as a harmonic oscillator with respect to their separation  $r$ . (B) Angles between neighbouring covalently bonded atoms are also approximated as a harmonic oscillator with respect to the angle  $\theta$ . In some forcefields such as CHARMM there is a correction term for these angular interactions known as Urey Bradley forces. This is calculated using the separation between the non-bonded atoms  $i-k$  in the triplet with the parameter  $r_{UB}$ . (C) The dihedral angle between four atoms is calculated by constructing two planes. Each plane is constructed to contain three of the four atoms in the set. One plane encompasses atoms  $i, j$  and  $k$  here colored in blue and the other plane contains the  $j, k$  and  $l$  atoms colored in red. The dihedral angle is then calculated by taking the angle between these two planes along the line they intersect, the line formed by the  $j-k$  bond. (D) The improper dihedral angles enforce the planarity of a molecular configuration. A plane is constructed to contain the  $i, j$  and  $k$  (blue) atoms and another plane is constructed to contain the  $j, k$  and  $l$  atoms (red). The improper angle is then calculated as the angle between these two planes.



**Figure 2.2: Comparison Between Potentials in Quantum and Classical Forcefields**

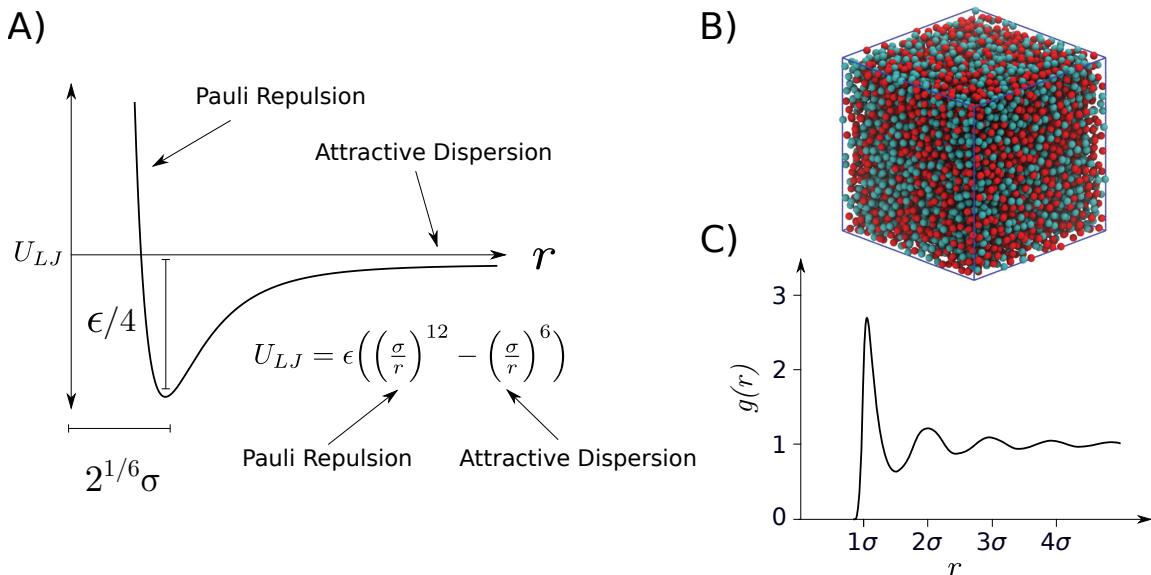
A) The Morse potential was formulated to approximate the potential the potential energy surface associated with the stretching of covalent bonds (blue). At low temperatures (the ground state,  $v = 0$ ) like those found in classical MD there is good agreement between the Morse potential and the harmonic oscillator (green). Credit Mark Somoza 2006 B) Here the potential of the dihedral angle between the atoms C1,C2,C3 and C4 in a butane molecule is calculated using two methods: Quantum Chemical calculations and approximations using the functional form in 2.14 [9]. Note how the appropriate choice of  $k_\varphi$ ,  $n$  and  $\delta$  have closely approximated the results the more accurate quantum mechanical calculations.

$$U_{non-bonded} = \underbrace{\sum_{i>j} \epsilon_{ij} \left( \left( \frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left( \frac{\sigma_{ij}}{r_{ij}} \right)^6 \right)}_{U_{Lennard-Jones}} - \underbrace{\sum_{i>j} \frac{q_i q_j}{r_{ij}}}_{U_{coulomb}} \quad (2.15)$$

The  $\sigma$  parameter denotes the location of the local minima in the Lennard-Jones potential. This is the optimum distance that two atoms will rest against each other in the absence of other effects. The  $\epsilon$  parameter denotes the depth of the potential well, or how stable the two atoms will be in the minimum energy configuration. This is very important for certain physical parameters such as osmotic pressure [8]

Conversely, the partial charges in a system have the greatest influence on the solvation energy.

By focussing on these two physical parameters we can isolate and improve the non-bonded parameters.



**Figure 2.3: The Lennard-Jones Potential**

A) The Lennard-Jones potential function has two regimes, the far region one dominated by attractive dispersion forces and the close region dominated by repulsion. In the case of atomic systems this is due to the Pauli exclusion principle. B) An example of a fluid modelled with Lennard-Jones particles [10]. C) The radial distribution function ( $g$ ) for a Lennard-Jones fluid [11]. Note that the peaks in the distribution are roughly  $1\sigma$  apart.

## The Lennard-Jones Potential

### 2.2.1 Philosophy of Different Molecular Mechanics forcefields.

At the time of writing, the four popular forcefields for the simulation of biomolecules are: AMBER, CHARMM, GROMOS and OPLS. Each of these have a slightly different philosophy in their formulation. They may be bottom up, as in the case of AMBER and CHARMM or top down, in the case of GROMOS and OPLS. Bottom up forcefields take the results from quantum *ab initio* calculations and approximate them with the functional form mentioned above. Conversely, top down forcefields take experimental measurable such as Osmotic pressure, solvation energy. This philosophy is closest to physics

### 2.2.2 Controlling the Temperature and Pressure in a Simulation

Living things are very sensitive to their external environment. Enzymes only work in a narrow range of temperatures and cells burst apart in the absence of pressure. As such, to correctly understand the events in biological bodies with simulations we not only need to correctly calculate the forces being exerted on every atom in their bodies but we must also make sure that the virtual environment in our simulations matches that found in the body or in the laboratory. Conceptually, we seek to approximate the environment as an open topped test-tube sitting in a pressure and temperature