# On Using Automated Algorithms to Parameterise Molecules for Molecular Dynamics Simulations
# and
# Investigating Suitable Ensembles for the Simulation of Naphthalimide Monolayers

Ivan Welsh

MASSEY UNIVERSITY
*Te Kunenga Ki Purehuroa*

**MASSEY UNIVERSITY**
TE KUNENGA KI PŪREHUROA
UNIVERSITY OF NEW ZEALAND

Institute of Natural and Mathematical Sciences

A thesis
submitted to Massey University in Albany, Auckland
in fulfilment of the requirements for the degree of
Doctor of Philosophy
in Chemistry.

Massey University Auckland
2017

# Abstract

Molecular dynamics simulations provide a means to investigate the spatial and temporal evolution of systems of molecules at atomic resolution. Force fields are used to describe the interactions between atoms contained within the system. A number of such force fields have been developed over the years, with a focus on force fields for use in simulations of biochemical systems, in particular, protein systems. This thesis is primarily focused on extending the range of systems that can be simulated through providing means for automated generation of force field parameters for large novel molecules.

One component of existing force fields that is generally poorly parameterised are the dihedral terms. In combination with the non-bonded terms, the dihedral terms are used to describe the rotational energy profile about bonds, and have a large influence on the conformational properties of a simulated system. A new method for the determination of dihedral parameters is developed, utilising high level quantum mechanical calculations. With the use of local elevation molecular dynamics simulations, this method is applied to the case of protein backbone dihedrals within the GROMOS force field.

When one desires to simulate the interaction of a novel molecule with some biochemical system, the novel molecule must be parameterised in a manner that is compatible with the force field used to describe the biochemical system. However, doing so is a slow, tedious, and error prone process, especially when the novel molecule is large. To combat this, a new algorithm, known as CherryPicker, was developed. CherryPicker is a graph based algorithm which enables rapid parameterisation of large molecules through fragment comparison with a library of previously parameterised small molecules. The algorithm design is discussed and tested on a few simple test cases in part II.

Part III steps away from the parameterisation focus of this thesis and looks at the simulation of naphthalimide monolayers. Naphthalimides have applications in sensing environments as they have absorption and fluorescence emission spectra lying within the UV and visible regions of light. With a long chain alkane substituted at the N-imide site, they become amphiphilic and can form monolayers on the surface of water, and can be transferred to a solid substrate when at a desired compression level. Molecular dynamics simulations can be used to provide insight into the formation of compressed monolayer phase. Here, the effect of different ensembles, namely NVT, NPT, and N$\gamma$T are investigated for use in simulating a naphthalimide monolayer.

# Acknowledgements

# Preface

With a preceding introductory chapter, this thesis is divided into three distinct parts. Part I focuses on the SpinningTop program, which is a program developed for determining dihedral parameters. A brief background to the reasons that such a method is required is given in chapter 2. Chapter 3 outlines the theory and implementation of the fitting method, and investigates some of the considerations that need to be made. Chapter 4 details the methods used to translate the developed fitting method to the case of protein backbone dihedral terms within the GROMOS force field, and chapter 5 discusses the results obtained in this proof of principle work.

Part II of the thesis focusses on the CherryPicker algorithm, which is a new algorithm developed to enable easy parameterisation of large biochemical molecules compatible with the GROMOS force field. As the algorithm is based on the concept of molecular fragmentation, a brief introduction to the state of computational molecular fragmentation is given in chapter 6. The design and mathematical background of the algorithm is presented in chapter 7, before a small amount of proof-of-concept testing is undertaken in chapter 8. As part of the CherryPicker algorithm development, a novel means to automatically determine bond order and formal charges of molecules was developed. This is presented in chapter 9.

Finally, part III presents work undertaken in the determination of suitable ensembles for the simulation of naphthalimide monolayers.

A large amount of code was developed as part of the work for this thesis. This code is available on request to j.allison@massey.ac.nz. The code will be provided as is, with no documentation on system requirements, installation, or usage.

# Contents

# List of Figures

# List of Tables

# List of Abbreviations

**AA** all atom.

**ATB** Automated Topology Builder.

**BSSE** basis set superposition error.

**CIP** Cahn–Ingold–Prelog priority rules.

**CSD** Cambridge Structural Database.

**DAG** directed acyclic graph.

**FPT** fixed parameter tractable.

**GAFF** Generalised AMBER force field.

**NBO** Natural Bond Orbital.

**PDB** Protein Data Bank.

**QM** quantum mechanical.

**QMDFF** Quantum Mechanically Derived Force Field.

**RMSD** root-mean-square deviation.

**UA** united atom.

**UFF** Universal Force Field.