# Development of a generic, structural bioinformatics information management system and its application to variation in foot-and-mouth disease virus proteins

by

Tjaart Andries Petrus de Beer

Submitted in partial fulfilment of requirements for the degree Philosophiae Doctor
(Bioinformatics)
in the Faculty of Natural and Agricultural Sciences
Bioinformatics and Computational Biology Unit
Department of Biochemistry
University of Pretoria
Pretoria
November 2008

# Declaration

I, Tjaart Andries Petrus de Beer, declare that the thesis/dissertation, which I hereby submit for the degree Philosophiae Doctor at the University of Pretoria, is my own work and not previously been submitted by me for a degree at this or any other tertiary institution.

SIGNATURE ......................................................... DATE ........................

# Acknowledgments

I want to thank the following people:

- My parents who supported me through all my studies.
- My supervisors for all their guidance, support and help during the last few years.
- All the various funding agencies who made it possible for me to study.
- All my friends and special people for their valuable support.
- My fellow students at the BCBU over the years.
- Irene, ti ringrazio per essermi stata accanto, sostenendomi sempre. Te ne sono davvero grato.

# Summary

Structural biology forms the basis of all functions in an organism from how enzymes work to how a cell is assembled. *In silico* structural biology has been a rather isolated domain due to the perceived difficulty of working with the tools. This work focused on constructing a web-based Functional Genomics Information Management System (FunGIMS) that will provide biologists access to the most commonly used structural biology tools without the need to learn program or operating specific syntax. The system was designed using a Model-View-Controller architecture which is easy to maintain and expand. It is Python-based with various other technologies incorporated. The specific focus of this work was the Structural module which allows a user to work with protein structures. The database behind the system is based on a modified version of the Macromolecular Structure Database from the EBI. The Structural module provides functionality to explore protein structures at each level of complexity through an easy-to-use interface. The module also provides some analysis tools which allows the user to identify features on a protein sequence as well as to identify unknown protein sequences. Another vital functionality allows the users to build protein models. The user can choose between building models online or downloading a generated script. Similar script generation utilities are provided for mutation modelling and molecular dynamics. A search functionality was also provided which allows the user to search for a keyword in the database. The system was used on three examples in Foot-and-Mouth Disease Virus (FMDV). In the first case, several FMDV proteomes were reannotated and compared to elucidate any functional differences between them. The second case involved the modelling of two FMDV proteins involved in replication, 3C and 3D. Variation between the several different strains were mapped to the structures to understand how variation affects enzymes structure. The last example involved capsid protein stability differences between two subtypes. Models

were built and molecular dynamics simulations were run to determine at which protein structure level stability was influenced by the differences between the subtypes. This work provides an important introductory tool for biologists to structural biology.

# Contents

# List of Abbreviations

| | |
|---|---|
| Å | Angstrom |
| aa/AA | Amino Acid |
| A | Alanine |
| ANSI | American National Standards Institute |
| C | Cysteine |
| CHARMM | Chemistry at HARvard Macromolecular Mechanics |
| CG | Conjugate Gradient |
| D | Aspartic acid |
| DNA | Deoxyribonucleic Acid |
| E | Glutamic acid |
| EBI | European Bioinformatics Institute |
| EC | Enzyme Commission |
| Ec | *Escherichia coli* |
| EM | Electron Microscopy |
| EST | Expressed Sequence Tag |
| F | Phenylalanine |
| FMDV | Foot and Mouth Disease Virus |
| FuGE | Functional Genomics Experiment |
| FunGIMS | Functional Genomics Information Management System |
| G | Glycine |
| GB | Gigabytes |
| H | Histidine |
| HAV | Hepatitis A Virus |
| HRV | Human Rhino Virus |

| | |
|---|---|
| HS | Heparan Sulfate |
| HMM | Hidden Markov Model |
| I | Isoleucine |
| ISO | International Standards Organization |
| K | Kelvin |
| K | Lysine |
| kb | kilobases |
| kD | kilo Dalton |
| L | Leucine |
| M | Methionine |
| MB | Megabytes |
| MSD | Macromolecular Structure Database |
| MVC | Model View Controller |
| N | Asparagine |
| NCBI | National Center for Biotechnology Information |
| ns | nanosecond |
| P | Proline |
| PDB | Protein Data Bank |
| Pfam | Protein families database |
| ps | picosecond |
| Q | Glutamine |
| R | Arginine |
| S | Serine |
| SD | Steepest Descent |
| sid | System Identifier |
| SQL | Structure Query Language |
| T | Threonine |
| TMHMM | Trans-Membrane Hidden Markov Model |
| V | Valine |
| W | Tryptophane |

XML        eXtensible Markup Language

Y          Tyrosine

# List of Figures

# List of Tables