

Accelerating the Interplay Between Theory and Experiment in Protein Design

Thesis by

Alex Nisthal

In Partial Fulfillment of the Requirements for the Degree of

Doctor of Philosophy

California Institute of Technology

Pasadena, California

2012

(Defended January 31, 2012)

© 2012

Alex Nisthal

All Rights Reserved

Acknowledgements

I'd first like to thank Steve Mayo for all of the support and advice he has given me over the years, including the opportunity to work on large expensive equipment with lots of bells and whistles.

I would also like to thank the rest of my thesis committee (Frances, Doug, and Dave) and Jost Vielmetter at the Protein Expression Center for providing direction when I needed it the most.

As the bridge between two distinct eras of the Mayo Lab, I've witnessed a complete overhaul of students and postdocs. From my time as a learner, I would especially like to thank Oscar Alvizo, Christina Vizcarra, Marie Ary, Rhonda Digiusto, Corey Wilson, Roberto Chica, Heidi Privett, Ben Allen, and Jennifer Keeffe for their time and geniality in teaching me the ways of the lab. Now, as the master, I'd like to thank Matt Moore, Kurt Mou, Alexandria Berry, Toni Lee, Tim Wannier, Bernardo Sosa Padilla Araujo, Mohsen Chitsaz, Gene Kym, Ernest Lee, Grace Lee, Samy Hamdouche, Emzo de los Santos, and Seth Lieblich for making the lab a fun, collaborative, and politically incorrect place to work.

Lastly, I would like to thank my friends for keeping me sane and my family for letting me be whatever I wanted to be.

Abstract

Protein engineering techniques such as directed evolution and structure-based design aim to improve the properties of natural proteins. The next step, the *de novo* insertion of function into previously inert protein scaffolds, is the lofty promise of computational protein design. In order to achieve this goal reliably and efficiently, computational methods can be iteratively improved by cycling between theory and experiment.

Efforts to both accelerate the rate and broaden the information exchanged within protein design cycles form the core of this thesis. Improvements in the throughput of experimental stability determination allowed the thorough assessment of new multi-state and library design tools. Intending to alleviate the fixed backbone, single native state design approximation, the study found constrained molecular dynamics ensembles useful for core repacking applications. The subsequent development of automated liquid handling protocols for common molecular biology techniques brings design experiments to new levels of sample throughput. This technology facilitated the creation of a stability database encompassing every single mutant in a small protein domain. Although constructed to facilitate future computational training efforts, we answer a multitude of questions pertaining to mutational outcomes, distributions, positional sensitivity, tolerance, and additivity in the context of a protein domain.

By expanding the constraints of experimental molecular biology, this work opens up new possibilities in the efforts to train and assay new computational methodologies for protein engineering applications.

TABLE OF CONTENTS

Acknowledgements		iii
Abstract		iv
Table of Contents		v
Figures and Tables		vi
Chapters		
Chapter 1	<i>Introduction</i>	1
Chapter 2	<i>Experimental library screening demonstrates the successful application of computational protein design to large structural ensembles</i>	10
Chapter 3	<i>Automated techniques for the complete site-directed mutagenesis and stability analysis of protein domains</i>	47
Chapter 4	<i>Stability analysis of the complete single mutant library of a protein domain</i>	80
Appendix		
Appendix	<i>High-throughput and automation methods</i>	126

Figures and Tables

Figure 2-1. The core residues of G β 1 designed in this study	35
Figure 2-2. General scheme used to design combinatorial mutation libraries based on computational protein design calculations	36
Table 2-1. Combinatorial libraries designed from sources of structural information	37
Table 2-2. Library coverage	38
Figure 2-3. Fraction-unfolded curves derived from the stability determination of experimental libraries	39
Figure 2-4. Library mutants sorted by experimental stability	40
Table 2-3. Combinatorial libraries designed from the top 16 energy-ranked structures based on two different energy functions	41
Figure 2-5. Library member energies	42
Figure 2-6. Correlation between simulation energy and experimental stability for the cMD-128 library	43
Figure 2-7. Microtiter plate-based stability assay controls	44-5
Figure 3-1. The automated site-directed mutagenesis pipeline	71
Figure 3-2. Visualization of a 96 well plate of SDM products	72
Figure 3-3. Variant construction timeline	73
Figure 3-4. Potential protein unfolding curves	74
Figure 3-5. Precision among experimental measures of protein stability	75
Figure 3-6. Dataset accuracy from literature comparisons	76
Figure 3-7. Point mutant amino acid distributions	77
Figure 4-1. Single mutant stability distribution for the G β 1 domain	108
Table 4-1. Gaussian fitting parameters for the mutational distributions	109
Figure 4-2. Gaussian fits of the G β 1 mutational distribution	110
Figure 4-3. Single mutant stability landscape for the GB1 domain	111
Figure 4-4. Single mutant stability distributions by RESCLASS	112
Figure 4-5. Packing density is linearly correlated with $\Delta\Delta G$ averaged by position	113
Figure 4-6. Amino acid scanning mutagenesis	114
Figure 4-7. Stability distribution of G β 1 by mutant amino acid	115
Figure 4-8. Calculated stability distributions by mutant amino acid	116
Table 4-2. Bioinformatics statistics for selected proteins	117
Table 4-3. Comparing the average $\Delta\Delta G$ of hydrophobic mutations by OSP	118
Table 4-4. Algorithm performance by linear correlation	119
Table 4-5. Algorithm performance by fraction correct	120
Figure 4-9. Complex additivity in core and surface mutation libraries	121
Table 4-6. Identity and stability of additive variants	122