Protein Interactions and Regulation of STIM and ORAI Genes

Dissertation

Zur Erlangung des Grades des Doktors der Naturwissenshaften der Naturwissenschaftlich-Technischen Fakultat III der Universität des Saarlandes

> von Ruzianisra Binti Mohamed

> > Saarbrücken 2017

Tag des Kolloquiums: Dekan: Perichterstatter:

Vorsitz: Akad. Mitarbeiter: 22.02.2017 Prof. Dr. Guido Kickelbick Prof. Dr. Volkhard Helms Prof. Dr. Richard Zimmermann Prof. Dr. Uli Müller Dr. Jessica Hoppstädter

ABSTRACT

Protein interactions play major roles in many biological processes. This thesis is composed of three projects. Using five datasets, we explored the characteristics and composition of overlapping protein-protein (PP) and protein-ligand (PL) interfaces. Overall, characteristics of PP contacts and overlapping PL contacts are highly similar. Second study was designed to identify transcription factor binding site motifs in promoter regions of STIM and ORAI genes, to gain knowledge of their regulation and relation with breast cancer. Our findings form an important basis of predictive interactions between transcription factors targeting STIM and ORAI genes and underline roles of these genes in breast cancer. Thirdly, we evaluated the performance of seven protein prediction tools on a dataset of protein-ligand complexes. Although the tools predicted pockets of various sizes and shapes, we found comparable performance amongst the predictions of five tools. We trained a random forest model to output a list of suitable tools for a given protein structure. This classifier should be useful for prioritizing the tools to be used for unknown proteins.

ZUSAMMENFASSUNG

Wechselwirkungen zwischen Proteinen spielen in biologischen Prozessen eine wesentliche Rolle. Diese Dissertation ist in drei Projekte aufgegliedert. Unter Zuhilfenahme von fünf Datensätzen aus drei unterschiedlichen Datenbanken wurden die Eigenschaften und die Zusammensetzung von überlappenden Protein-Protein (PP) und Protein-Ligand (PL) Bindestellen untersucht. Statistisch gesehen sind sich die Eigenschaften von PP Kontakten und überlappenden PL Kontakten sehr ähnlich. Die zweite Studie diente dazu Motive von Transkriptionsfaktor-Bindestellen in bestimmten Promotor-Regionen der STIM und ORAI Gene zu identifizieren als auch Erkenntnisse über deren Regulation und Zusammenhang mit Brustkrebs zu gewinnen. Unsere Ergebnisse stellen eine wichtige Grundlage für die Vorhersage der Wechselwirkungen zwischen Transkriptionsfaktoren, welche an die beiden Gene STIM und ORAI binden, dar. Darüber hinaus konnten wir die Rolle dieser Gene in Zusammenhang mit Brustkrebs herausstellen. Im dritten Teil werteten wir die Performance von sieben verschiedenen Tools anhand eines Datensatzes von PL Komplexen aus. Obwohl die Tools von der Größe und Form her unterschiedliche Bindetaschen vorhersagten, konnten wir dennoch ein vergleichbares Verhalten zwischen fünf Tools feststellen. Wir trainierten ein Random-Forest Modell, welches, gegeben eine Proteinstruktur, eine Reihe von geeigneten Tools vorhersagt. Dieses Modell kann dazu dienen Vorhersagemodelle anhand unbekannter Proteine zu priorisieren.

ACKNOWLEDGEMENTS

Throughout my PhD study and research, it has been a pleasure working with my great supervisor and colleagues, this dissertation owe much to their sincere help and encouragement. First, I would like to thank my PhD supervisor Prof. Volkhard Helms, for his support and encouragement throughout the whole PhD work. His rational and rigorous approach in doing scientific research has been significantly influential and greatly helpful for this dissertation.

Secondly, I had a great fun and a good collaboration from my colleagues and their contributions should not be underestimated. I thank my collaborators Jennifer Degac, Zhao Yuan, Rahmad Akhbar, Riccha Sethi, and Dr. Mohamed Hamed. A big thank you also to Maryam Nazeriah, Thorsten Will, Kerstin Reuter, Dr. Michael Hutter, Dr. Ruslan Akulenko, Dr. Siba Shanak, Tran Vu Ha, Dr. Ahmad Bargash, Dr. Ozlem Ulucan, Daria Gadar, Dr. Christian Spaniol, Shimaa Khaled, Jan Riehm, Duy Nguyen and Hema as they gave me a lot of help and suggestions. I also appreciate the secretary of my supervisor Kerstin Pudelek for her big help in administration works.

I am also grateful to Ministry of Education of Malaysia for providing me a SLAI/SLAB scholarship for Phd. And big thanks to MARA University of Technology (UiTM) and all staffs of Faculty of Pharmacy, UiTM Puncak Alam for supporting me to study abroad.

Finally, I would like to mention my father, my siblings, nieces, nephews, best friends, my fiancé and the biggest appreciations to my late Bonda and my brother. I am very grateful for the supports from the whole family and friends. Without their supports and prayers, I could not finish my PhD. They bring me a lot of happiness and give me encouragement to overcome the difficulties during my scientific research.

LIST OF FIGURES

| Figure 2.1 | Classification of different types of protein interactions |
|------------|--|
| Figure 2.2 | Pocket (orange sphere) identified by the protein pocket identification tool GHECOM that overlaps with the ligand dysiherbaine (red sticks) bound to human glutamate receptor, GluR5 (PDB ID: 3FV1) |
| Figure 2.3 | Example of formation of position weight matrices (PWM) and sequence motif. (A) Shown are eight known genomic binding sites in three <i>Saccharomyces cerevisiae</i> genes (HEM13, ANB1 and ROX1). (B) Frequencies matrix of nucleotides at each position. (C) Sequence logo used to visualize count at each position. (D) Sequence logo to represent the frequencies scaled relative using the information content at each position. (E) Energy normalized logo using relative entropy to adjust the GC content in <i>S.cerevisiae</i> |
| Figure 3.1 | Schematic illustration of a PL complex illustrating the interface (black border) and the remaining surface regions |
| Figure 3.2 | Flowchart summarizing the compilation of contacts between amino acids of the first protein (Pi1) and amino acids of the second protein (Pi2), atom contacts in PP and PL complexes, and the calculation of PR and IR |
| Figure 3.3 | Percentage frequencies and propensities of amino acid residues at protein interfaces of PP complexes from the ABC dataset. (A) Percentage frequencies of amino acid residues at protein interfaces. (B) Propensities of amino acid residues at protein interfaces |
| Figure 3.4 | Percentage frequencies and propensities of amino acid residues at protein interfaces of PP complexes from the PIBASE dataset. (A) Percentage frequencies of amino acid residues at protein interfaces. (B) Propensities of amino acid residues at protein interfaces |
| Figure 3.5 | Percentage frequencies and propensities of amino acid residues at protein interfaces of PL complexes from the ABC, PIBASE and Timbal datasets. (A) Percentage frequencies of amino acid residues at protein interfaces. (B) Propensities of amino acid residues at protein interfaces |
| Figure 3.6 | Amino acid pairing propensities (in log ₂ -format) for interfaces of PP complexes from the ABC dataset |
| Figure 3.7 | Amino acid pairing propensities (in log ₂ -format) for interfaces of PP complexes from the PIBASE dataset |
| Figure 4.1 | Molecular structure of STIM1 (A) STIM1 domain organization. (B) Cartoon depicting a possible model of the STIM1 monomer in the resting state |
| Figure 4.2 | Molecular structure of STIM2 (A) STIM2 domain organization. (B) Cartoon depicting a possible model of the STIM2 monomer in the resting state |

| Figure 4.3 | Molecular structures of ORAI1, ORAI2, and ORAI3 (A) Domain organization of ORAI genes. (B) Cartoon depicting a possible model of ORAI in the resting state and sequence alignments |
|-------------|--|
| Figure 4.4 | Flowchart summarizing the workflow of the analysis46 |
| Figure 4.5 | Dynamic graphical representations of 48 transcription factors targeting STIM and ORAI genes. The figure drawn using Cytoscape |
| Figure 4.6 | Schematic illustrations of transcription factors binding motifs in the promoters of the genes STIM1, STIM2, ORAI1, ORAI2 and ORAI3. (A) Transcription factor binding motifs in the promoters of genes STIM1 and STIM2. (B) Transcription factors binding motifs in the promoters of genes ORAI1, ORAI2 and ORAI350 |
| Figure 4.7 | Network of 14 transcription factors targeting STIM and ORAI genes. Evidence view of the STRING database output depicting the transcription factors targeting STIM and ORAI genes obtained from <u>http://string-db.org/</u> |
| Figure 4.8 | Plot of principal component analysis of 113 normal samples and 1102 tumor samples |
| Figure 4.9 | Topological overlap matrix (TOM) heatmap corresponding to the two co- expression modules in normal samples. Each row and column of the heatmap represents a single gene. Red indicates high levels of co-expression genes. The dendograms on the upper and left sides show the hierarchical clustering tree of genes |
| Figure 4.10 | Topological overlap matrix (TOM) heatmap corresponding to the three co- expression modules in tumor samples. Each row and column of the heatmap represents a single gene. Red indicates high levels of co-expression genes. The dendograms on the upper and left sides show the hierarchical clustering tree of genes |
| Figure 4.11 | Differential Interactions Networks (A) Differential interactions network of normal-tumor. (B) Differential interactions network of tumor-normal. Red edges indicate interactions that found in differential interactions and black edges indicate the interactions which were not found in differential interactions. Blue nodes indicate the genes found in the blue module, brown (blue module), and turquoise (turquoise module), respectively |
| Figure 5.1 | Two dimensional illustrations of three possible ways to define the pocket volume of a pacman-shape surface cavity |
| Figure 5.2 | Pockets (blue) identified by DEPTH (lining residues) (top panel left), GHECOM (top panel middle), Fpocket (top panel right), DoGSiteScorer (middle panel left), PocketPicker (middle panel middle), IsoMif (middle panel right), and ProACT2 (bottom panel) that overlap with the ligand HNT (red sticks) bound to human phenylethanolamine N-methyltransferase, PNMT (PDB ID: 2G70). The figures were generated using PyMOL Molecular Graphics System |

| Figure 5.3 | Distribution | of the | number | of | predicted | pockets | per | protein | for | Fpocket, |
|------------|--------------|----------|------------|------|-----------|---------|-----|---------|-----|----------|
| | DoGSiteSco | rer. and | IsoMif. re | espe | ectivelv | | | | | 72 |
| | | - , | , | | , | | | | | |

| Figure 5.4 | Distribution | of the number of | of pocket lini | ing residues pe | r pocket | (top five) for |
|------------|--------------|------------------|----------------|-----------------|----------|----------------|
| | Fpocket, | DoGSiteScorer, | IsoMif, | GHECOM, | and | PocketPicker, |
| | respectively | | | | | 73 |

LIST OF TABLES

| Table 2.1 | Classification of different types of protein interactions |
|-----------|--|
| Table 2.2 | Pocket (orange sphere) identified by the protein pocket identification tool GHECOM that overlaps with the ligand dysiherbaine (red sticks) bound to human glutamate receptor, GluR5 (PDB ID: 3FV1) |
| Table 2.3 | Example of formation of position weight matrices (PWM) and sequence motif. (A) Shown are eight known genomic binding sites in three <i>Saccharomyces cerevisiae</i> genes (HEM13, ANB1 and ROX1). (B) Frequencies matrix of nucleotides at each position. (C) Sequence logo used to visualize count at each osition. (D) Sequence logo to represent the frequencies scaled relative using the information content at each position. (E) Energy normalized logo using relative entropy to adjust the GC content in <i>S. cerevisiae</i> |
| Table 3.1 | Schematic illustration of a PL complex illustrating the interface (black border) and the remaining surface regions |
| Table 3.2 | Flowchart summarizing the compilation of contacts between amino acids of the first protein (Pi1) and amino acids of the second protein (Pi2), atom contacts in PP and PL complexes, and the calculation of PR and IR |
| Table 3.3 | Percentage frequencies and propensities of amino acid residues at protein interfaces of PP complexes from the ABC dataset. (A) Percentage frequencies of amino acid residues at protein interfaces. (B) Propensities of amino acid residues at protein interfaces |
| Table 3.4 | Percentage frequencies and propensities of amino acid residues at protein interfaces of PP complexes from the PIBASE dataset. (A) Percentage frequencies of amino acid residues at protein interfaces. (B) Propensities of amino acid residues at protein interfaces |
| Table 3.5 | Percentage frequencies and propensities of amino acid residues at protein interfaces of PL complexes from the ABC, PIBASE and Timbal datasets. (A) Percentage frequencies of amino acid residues at protein interfaces. (B) Propensities of amino acid residues at protein interfaces |
| Table 3.6 | Amino acid pairing propensities (in log ₂ -format) for interfaces of PP complexes from the ABC dataset |
| Table 3.7 | Amino acid pairing propensities (in log ₂ -format) for interfaces of PP complexes from the PIBASE dataset |
| Table 4.1 | Molecular structure of STIM1 (A) STIM1 domain organization. (B) Cartoon depicting a possible model of the STIM1 monomer in the resting state |
| Table 4.2 | Molecular structure of STIM2 (A) STIM2 domain organization. (B) Cartoon depicting a possible model of the STIM2 monomer in the resting state |

| Table 4.3 | Molecular structures of ORAI1, ORAI2, and ORAI3 (A) Domain organization of ORAI genes. (B) Cartoon depicting a possible model of ORAI in the resting state and sequence alignments |
|------------|--|
| Table 4.4 | Flowchart summarizing the workflow of the analysis46 |
| Table 4.5 | Dynamic graphical representations of 48 transcription factors targeting STIM and ORAI genes. The figure drawn using Cytoscape |
| Table 4.6 | Schematic illustrations of transcription factors binding motifs in the promoters of the genes STIM1, STIM2, ORAI1, ORAI2 and ORAI3. (A) Transcription factor binding motifs in the promoters of genes STIM1 and STIM2. (B) Transcription factors binding motifs in the promoters of genes ORAI1, ORAI2 and ORAI350 |
| Table 4.7 | Network of 14 transcription factors targeting STIM and ORAI genes. Evidence view of the STRING database output depicting the transcription factors targeting STIM and ORAI genes obtained from <u>http://string-db.org/</u> |
| Table 4.8 | Plot of principal component analysis of 113 normal samples and 1102 tumor samples |
| Table 4.9 | Topological overlap matrix (TOM) heatmap corresponding to the two co- expression modules in normal samples. Each row and column of the heatmap represents a single gene. Red indicates high levels of co-expression genes. The dendograms on the upper and left sides show the hierarchical clustering tree of genes |
| Table 4.10 | Topological overlap matrix (TOM) heatmap corresponding to the three co- expression modules in tumor samples. Each row and column of the heatmap represents a single gene. Red indicates high levels of co-expression genes. The dendograms on the upper and left sides show the hierarchical clustering tree of genes |
| Table 4.11 | Differential Interactions Networks (A) Differential interactions network of normal-tumor. (B) Differential interactions network of tumor-normal. Red edges indicate interactions that found in differential interactions and black edges indicate the interactions which were not found in differential interactions. Blue nodes indicate the genes found in the blue module, brown (blue module), and turquoise (turquoise module), respectively |
| Table 5.1 | Two dimensional illustrations of three possible ways to define the pocket volume of a pacman-shape surface cavity |
| Table 5.2 | Pockets (blue) identified by DEPTH (lining residues) (top panel left), GHECOM (top panel middle), Fpocket (top panel right), DoGSiteScorer (middle panel left), PocketPicker (middle panel middle), IsoMif (middle panel right), and ProACT2 (bottom panel) that overlap with the ligand HNT (red sticks) bound to human phenylethanolamine N-methyltransferase, PNMT (PDB ID: 2G70). The figures were generated using PyMOL Molecular Graphics System |

| Table 5.3 | Distribution of the number of predicted pockets per protein for Fpocket, DoGSiteScorer, and IsoMif, respectively72 |
|-----------|---|
| Table 5.4 | Distribution of the number of pocket lining residues per pocket (top five) for Fpocket, DoGSiteScorer, IsoMif, GHECOM, and PocketPicker, respectively |
| Table 5.5 | Cumulative density of Matthews correlation coefficient (MCC), precision (pre), recall, and residues overlap (overlap)74 |

ABREVIATIONS

| 3D | Three dimensional |
|-----------|--|
| Å | Angstrom |
| aa | Amino acid |
| bp | Base pair |
| BRCA | Breast invasive carcinoma |
| Ca^{2+} | Calcium |
| ChIP | Chromatin immuniprecipitation |
| ChIP-chip | Chromatin immuniprecipitation microarrays |
| CRAC | Calcium Release Activated Calcium |
| Da | Dalton |
| DESeq2 | Differential Expression analysis for Sequence count data version 2 |
| DNA | Deoxyribonucleic acid |
| edgeR | Empirical analysis of digital gene expression in R |
| ENCODE | Encyclopedia of DNA Elements |
| ER | Endoplasmic reticulum |
| FDR | False Discovery Rate |
| FN | False negative |
| FP | False positive |
| MD | Molecular Dynamics |
| MCC | Matthews coefficient correlation |
| ORAI1 | Calcium release-activated calcium channel protein 1 |
| ORAI2 | Calcium release-activated calcium channel protein 2 |
| ORAI3 | Calcium release-activated calcium channel protein 3 |
| PDB | Protein Data Bank |
| PL | Protein-ligand |
| PP | Protein-protein |
| PWM | Position weight matrices |
| SASA | Solvent Accessible Surface Area |
| SOCE | Store-Operated Ca ²⁺ entry |
| STIM | Stromal Interaction Molecules |
| STIM1 | Stromal Interaction Molecule 1 |
| STIM2 | Stromal Interaction Molecule 2 |
| SVM | Support Vector Machine |
| TCGA | The Cancer Genome Atlas |
| TF | Transcription factor |
| TFBS | Transcription factor binding site |
| TM | Transmembrane |
| TN | True negative |
| TOM | Topological overlap matrix |
| TP | True positive |
| TSS | Transcription start site |
| WGCNA | Weight Gene Co-expression Network Analysis |
| | |

TABLE OF CONTENTS

| Abstract | i |
|---|-------|
| Zusammenfassung | ii |
| Acknowledgements | . iii |
| List of Figures | . iv |
| List of Tables | vii |
| Abreviations | x |
| Table of Contents | . xi |
| | |
| Chapter 1: Introduction | 1 |
| 1.1 Introduction | 1 |
| 1.2 Motivation | 1 |
| 1.3 Contributions | 2 |
| 1.4 Thesis Organization | 3 |
| | |
| Chapter 2: Biological Background and Computational Methods | 4 |
| 2.1 The Nature of Protein | 4 |
| 2.2 Protein-Protein Interactions, Protein-Ligand Interactions, and Protein Interfaces | 4 |
| 2.3 Types of Protein Interactions | 5 |
| 2.4 The Protein Data Bank (PDB) | 6 |
| 2.5 Protein Pockets | 7 |
| 2.6 Computational Methods for Identification of Protein Pockets | 8 |
| 2.6.1 Geometry based Methods | 8 |
| 2.6.1.1 Grid System Scanning | 9 |
| 2.6.1.2 Probe Sphere Filling Methods | 9 |
| 2.6.1.3 Alpha Shape Methods | 10 |
| 2.6.2 Energy based Methods | 10 |
| 2.6.3 Evolution based Methods | 10 |
| 2.6.4 Blind Docking and Molecular Dynamic Methods | 11 |
| 2.6.5 Combined Approaches | 11 |
| 2.7 Transcription Factors, Promoters, Motifs, and Position Weight Matrix (PWM) | .12 |
| 2.7.1 Transcription Factors and Promoters | 12 |
| 2.7.2 Motifs and position weight matrix (PWM) | 12 |
| 2.8 Gene Expression, Differential Gene Expression, and Co-expression Analysis | 13 |
| 2.8.1 Gene Expression | 13 |
| 2.8.2 Differential Gene Expression Analysis | 14 |
| 2.8.2.1 DESeq. DESeq2. and edgeR Packages | 14 |
| 2.8.3 Co-expression Analysis | 15 |
| | |
| Chapter 3: Composition of Overlapping Protein-Protein and Protein-Ligand Interfaces | 17 |
| Abstract | 17 |
| 3.1 Introduction | 18 |

| 3.2 Material and Methods | 19 |
|---|----|
| 3.2.1 Datasets | 19 |
| 3.2.2 Surface and Interface Residues | 20 |
| 3.2.3 Classification of the Amino Acids | |
| 3.2.4 Interface Residue Propensities | |
| 3.2.5 Surface and Interface Residues | |
| 3.2.6 Atom Contacts in Protein-protein and Protein-ligand Complexes | |
| 3.2.7 Calculation of Polarity Ratio and Interface Atom Ratio | 22 |
| 3.3 Results and Discussion | 22 |
| 3.3.1 Amino Acid Composition and Protein Interfaces Propensity | |
| 3.3.2 Amino Acid Contacts | |
| 3.3.3 Atomic Contacts in Protein-protein and Protein-ligand Complexes | 32 |
| 3.3.4 Polarity Ratio and Interface Atom Ratio | 37 |
| 3.4 Conclusions | |

| Chapter 4: STIM and ORAI Genes, Interactions with Transcription Factors, Differential |
|---|
| Gene Expression and Co-expression Analysis on Breast Invasive Carcinoma |
| Dataset |
| Abstract |
| 4.1 Introduction |
| 4.2 Material and Methods |
| 4.2.1 Transcription Factors Targeting STIM and ORAI Genes |
| 4.2.2 Transcription Factors Binding Models 43 |
| 4.2.3 Sequence of Promoter Regions 44 |
| 4.2.4 Motif Over-representation Analysis 44 |
| 4.2.5Transcription Factors Predicted by CheA Mapped with the FIMO Results44 |
| 4.2.6 Prediction of Physical Interactions of Protein-protein Complexes |
| 4.2.7 TCGA BRCA Dataset 45 |
| 4.2.8 Computing Differential Analysis 45 |
| 4.2.9 Weight Gene Co-expression Network Analysis 45 |
| 4.2.10 Flowchart |
| 4.3 Results and Discussion |
| 4.3.1 Transcription Factors Targeting STIM and ORAI Genes |
| 4.3.2 Promoter Sequences of STIM and ORAI Genes from the EPDnew Database 48 |
| 4.3.3 HOCOMOCO and FIMO Results |
| 4.3.4 Physical Protein-protein Interactions |
| 4.3.5 Differential Gene Expression Analysis |
| 4.3.6 Results of Gene Co-expression Analysis |
| 4.3.6.1 Normal samples |
| 4.3.6.2 Tumor samples |
| 4.3.7 Differential Interaction Networks |
| 4.3.8 Regulation of STIM and ORAI Genes on Normal and Tumor Breast Cancer |
| Tissue |
| 4.4 Conclusions |

| Chapter 5: Evaluation of Protein Pocket Identification Tools on Protein-I | Ligand Complexes |
|---|------------------|
| Abstract | |
| 5.1 Introduction | |
| 5.2 Material and Methods | 68 |
| 5.2.1 Dataset | 68 |
| 5.2.2 Tools | 68 |
| 5.2.3 Binding Site Definition | |
| 5.2.4 Model Evaluations | |
| 5.2.5 MCC | |
| 5.2.6 Precision | |
| 5.2.7 Recall | |
| 5.2.8 Overlap | |
| 5.2.9 Correlation between Features and Quality Metrics | |
| 5.2.10 TSore | |
| 5.2.11 TRatio | |
| 5.2.12 Classifiers | |
| 5.3 Results | 71 |
| 5.3.1 Pocket Distributions | 71 |
| 5.3.2 Size Distributions | |
| 5.3.3 Prediction versus Reality | |
| 5.3.4 Features and Quality Metrics | 74 |
| 5.3.5 Tool Score | |
| 5.3.6 Identifying an Optimal Tool for a Protein | |
| 5.4 Discussion | |
| 5.5 Conclusions | |
| Chapter 6: Conclusions and Future Works | |
| 7: References | |
| 8: Supplementary Information | |
| 8.1 Supplementary Information for Chapter 3 | |
| 8.2 Supplementary Information for Chapter 4 | |
| 8.3 Supplementary Information for Chapter 5 | |

Chapter 1 Introduction

1.1 Introduction

Protein interactions play a major role in many biological processes. Numerous important applications benefit from the identification of protein-protein and protein-ligand interactions such as drug design, protein mimetics engineering, elucidation of molecular pathways, and understanding of disease mechanisms. Nowadays, a wide range of high-throughput experimental approaches are available for identification of protein-protein and protein-ligand interactions. For that reasons, many publicly accessible databases stores the high-quality information about protein-protein and protein-ligand complexes. Moreover, for almost three decades, gene expression was recognized to be mainly regulated at the transcriptional level and protein known as transcription factors functions in regulating the expression of a gene. Any changes of protein interactions can lead to mutations and diseases by affecting the functions of protein complexes or by affecting gene regulations.

1.2 Motivation

This thesis addresses three problems which related mostly to protein interactions and gene regulation. Nowadays, many efforts and studies have analyzed protein-protein and protein-ligand interactions, however a detailed of protein interfaces composition is remain elusive as targeting protein-protein interactions is challenging because usually no convenient natural substrates are available as starting point for small-molecule design. Hence, by using five datasets from three different databases, our aim was to explore the characteristics and composition of overlapping protein-protein and protein-ligand interfaces.

On the other hand, many studies have been done on STIM and ORAI genes as they play important roles in calcium signals and involved in store-operated Ca^{2+} entry (SOCE) in cells. STIM1 and STIM2 are needed for the development and functioning of various cell types such as lymphocytes, skeletal and smooth muscle myoblast, adipocytes and neurons [1,2]. In addition, previous studies stated that STIM and ORAI genes are associated with differentiation processes and linked to several diseases such as Alzheimers, Parkinsons diseases and cancer (e.g. breast cancer and prostate cancer) [3–6]. Today, cancer is a major public health problem worldwide which breast cancer is one of the most common and predominant cancer types that affects 1.3 million people and causes thousands of death annually [7,8]. In the United States of America, 30.4% of women are diagnosed with breast cancer yearly [7]. However, only few studies addressed the involvements of STIM and ORAI genes to breast cancer. In addition, the discovery of transcription factor binding sites (TFBS) motifs in specific locations in the promoter regions of STIM and ORAI genes is remain elusive. Our aim was to find the interacting transcription factors bound to the particular promoter regions of STIM1, STIM2, ORAI1, ORAI2, and ORAI3 and their regulation and relation with normal and breast invasive carcinoma (BRCA) samples.

Identification of binding pockets is often a prerequisite step for structure-based drug design. However, the characterization of these sites is a main challenge in computational biology. Defining the correct pockets on protein surfaces is not an easy task. Recently, many protein pocket identification tools have been developed by using different strategies and approaches. So far apparently, none of the existing approaches performed in a way which can set as a standard or that is widely accepted as benchmarking method to identify protein pockets accurately. Hence, it is great of interest to evaluate and compare the performances of these tools. By using seven different protein pocket identification tools and tested on a dataset of protein-ligand complexes, our aim was to evaluate and compare their performances and suggest a set of the best tool which can identify protein pockets accurately.

1.3 Contributions

Most of the results chapters of this thesis are based on manuscripts that either published or ready for submissions as follow:

Chapter 3: Ruzianisra Mohamed, Jennifer Degac, and Volkhard Helms. Composition of Overlapping Protein-Protein and Protein-Ligand Interfaces. PloS One. 2015 Oct 30;10(10):e0140965. doi: 10.1371/journal.pone.0140965.

Chapter 4: Ruzianisra Mohamed, Riccha Sethi, Mohamed Hamed, and Volkhard Helms. STIM and ORAI Genes, Interactions with Transcription Factors, Differential Gene Expression and Co-expression Analysis on Breast Invasive Carcinoma Dataset (In preparation for submission to peer reviewed journal). Chapter 5: Zhao Yuan, Rahmad Akbar, Volkhard Helms, and Ruzianisra Mohamed. Evaluation of Protein Pocket Identification Tools on Protein-Ligand Complexes (In preparation for submission to peer reviewed journal).

1.4 Thesis Organization

The structure of the thesis is as follows:

- Chapter 2 provides general introduction to biological background and computational methods such as proteins, protein-protein and protein-ligand interactions, protein pocket identification algorithms, transcription factors, promoters, gene expressions, differential gene expression analysis and co-expression analysis.
- Chapter 3 will discuss the topic on compositions of overlapping protein-protein and protein-ligand interfaces.
- Chapter 4 will show the study on STIM and ORAI genes, interactions with transcription factors and differential gene expression and co-expression analysis on a dataset of breast invasive carcinoma.
- Chapter 5 will discuss the topic of evaluation of seven protein pocket identification tools on a dataset of 167 protein-ligand complexes.
- Chapter 6 summarizes the results of the projects and provides conclusions with regard to the aims of the studies and contributions made and present the future works.

Chapter 2

Biological Background and Computational Methods

In this chapter, we present the background on the biological concepts and computational methods relevant to this thesis. We also briefly describe the publicly accessible biological databases, tools, and packages that were used in this thesis.

2.1 The Nature of Protein

Proteins rarely act by themselves and in most cases in order to function in biological systems, they work in groups which known as protein complexes. Protein complexes perform many important tasks within biological cells such as catalyzing metabolic reactions, replicating DNA, and transporting molecules from one place to another. Being workhorses that assist many biological processes, their detection is needed to increase our knowledge about cellular organizing and function. Generally, a protein consists of a sequence of 20 amino acids (also known as residues) which are linked together through peptide bonds, forming a polypeptide chain known as the primary structure. The primary structure forms secondary structural elements, for example alpha helices and beta sheets through hydrogen bonds, which interact to form the tertiary structure via protein folding. The tertiary structure may combine with another to construct a quaternary structure.

2.2 Protein-Protein Interactions, Protein-Ligand Interactions, and Protein Interfaces

As mentioned above, proteins play major roles in almost all biological functions. Most of the molecular processes are based on molecular machines which are composed of a large number of proteins and bind to each other through protein-protein interactions. The interactions of proteins are determined and usually mediated by their interfaces [9]. Protein interfaces are defined as binding sites of certain patches on each of the two protein's surface, which enables the interaction between two proteins. The characteristics of the protein-protein and protein-ligand interfaces are important to describe binding principles and provide clues about algorithms for binding site prediction. One of the most important aims in interface studies is to identify properties which may distinguish residues at binding sites from the rest of the protein surfaces. A few studies showed that interfaces are rather large, planar or well packed depending on the type of the interactions [10,11].

Due to the importance of protein interfaces in mediating protein-protein and proteinligand interactions, various studies have analyzed protein interfaces. The analysis of three dimensional (3D) protein structures revealed that protein interfaces are composed of buried cores which are surrounded by partially accessible rims [12,13]. Additionally, the size of the interface patches ranges between 400 and 1600Å² [12].

Interestingly, PPI interfaces contain small regions called "hot spots" that contribute the most to the total binding energy [10,14,15]. A powerful experimental approach named alanine scanning mutagenesis is used to identify hot spot regions by successively mutating each interface residue to alanine. Those residues are termed hot spots where mutation into alanine leads to a decrease in binding free energy by more than 2.0 kcal/mol [16]. Since hot spots by definition make up the largest contribution to the binding energy, the ability to predict hot spots is important to identify, analyze and lead to targets for drug binding sites. To circumvent the required time and costs and the experimental uncertainty whether protein mutants will properly express and purify, several computational methods have been designed and used to predict hot spots. They typically either use energy based methods to calculate the energetic contributions of residues [17,18], solvent accessible surface area (SASA) [19], and sequence conservation based approaches [20].

2.3 Types of Protein Interactions

Studies showed that many factors are influencing the classification of protein interactions into different types of interactions [10,21]. The most basic classification of interactions is based on the composition of the protein complex. For example a complex consisting of identical proteins is termed homo-oligomer (e.g. homodimer complex). Usually, homo-oligomers are very stable permanent protein structures such as in the case of homodimer (two proteins) complexes [10]. On the other hand, a complex made of non-identical proteins is defined as hetero-oligomers (e.g. heterodimer complex) [22]. Then, protein interactions can be further grouped into obligate or non-obligate depending on the lifetime of the interactions [22,23]. Obligate complexes only exist in the bound form. Components of non-obligate interactions are also stable in the unbound form. Based on their binding affinity, the interaction can be considered to be permanent or transient. Generally, all obligate interactions are permanent, however not all permanent interactions are obligate. Non-obligate and transient protein interactions usually have lower binding affinities. The transient interactions are subdivided into weak and strong transient. The strong transient ones include protein interactions that shift from an unbound or weakly bound state to a strongly bound state. This is frequently stimulated by an effector

molecule. Figure 2.1 shows the classification of protein interactions. The protein interaction types depend on the situation and cellular process. Therefore, it is importance to understand protein-protein interaction types and the effects they have on certain biological processes.



Figure 2.1 Classification of different types of protein interactions.

2.4 The Protein Data Bank (PDB)

The 3D coordinates of protein structures are generally deposited in The Protein Data Bank (PDB) [24], the globally recognized primary depository for experimentally determined atomistic structure of 3D biological macromolecules. The PDB was developed in 1971 at the Brookhaven National Laboratories containing a set of seven protein structures. Then, in 1998, the PDB moved under management of the Research Collaboratory for Structural Bioinformatics (RCSB) at the Rutgers University, New Jersey (http://www.rcsb.org/pdb/). Since early 1990's, the number of 3D structures deposited in the PDB has been increasing exponentially. As for 16 May 2016, the PDB stores 118748 structure of biological macromolecules of which more than 92% are proteins, 2% nucleic acid complexes, 4% protein-nucleic acid complexes and 0.02% other complexes. 99% of the structures were determined by X-ray crystallography (89%) and by Nuclear Magnetic Resonance (NMR) Spectroscopy (10%). 1022 structures were solved with electron microscopy, 95 hybrid and 196 by other methods. The PDB also stores 20955 ligands which interact with the proteins (protein-ligand complexes).

The PDB is the primary resource to study the diversity of protein-protein (PP) and protein-ligand (PL) interactions. Thus, some secondary databases have been developed to

assist the research on PP interactions such as the Biological General Repository for Interaction Datasets BioGRID (thebiogrid.org) [25,26] which stores data of small molecules modulating protein-protein **INTeraction** (MINT) complexes, the Molecular database (mint.bio.uniroma2.it/mint/Welcome.do) [27], the Biomolecular Interaction Network Database (BIND) (http://bind.ca) [28], the Database of Interacting Proteins (DIP: http://dip.doembi.ucla.edu) [29], the IntAct molecular interaction (IntAct) database (www.ebi.ac.uk/intact/) [30], the ABC database (http://service.bioinformatik.uni-saarland.de/ABCSquareWeb/) [31], database of structurally defined protein interfaces PIBASE and the named (http://pibase.janelia.org/pibase2010/queries.html) [32].

On the other hand, the high-quality databases of PL interactions can assist the field of structure-based drug design to develop the best computational tools. There are several publicly available databases such as the Timbal database (<u>http://mordred.bioc.cam.ac.uk/timbal</u>) [33], the Mother of All Database (MOAD) (<u>http://bindingmoad.org</u>) [34,35], the 2P2I database (<u>http://2p2idb.cnrs-mrs.fr</u>) [36], the PDBbind (<u>http://2p2idb.cnrs-mrs.fr</u>) [36], the PDBbind (<u>http://sw16.im.med.umich.edu/databases/pdbbind/index.jsp</u>) [37], the BindingDB database (<u>www.bindingdb.org</u>) [38,39], and Relibase (<u>http://relibase.ccdc.cam.ac.uk/account_utilities/login_form.php</u>) [40]. The ABC and PIBASE databases also stored data on PL interactions.

2.5 Protein Pockets

The accurate prediction of protein pockets from a 3D protein structure is an important issue for structure-based drug design [41,42] which can elucidate protein functions [43]. The basic principle of molecular interactions was proposed by Emil Fischer in 1894, who stated that ligand binding to proteins follows a "lock and key" mechanism [44]. These interactions often occur in particular sites on protein surfaces (binding sites). These binding sites can be distinguished from other parts of the protein surface by their unique characteristics as such the electrostatic properties and the size of a cavity on the protein surface [45]. Additionally, in a single protein there can be many pockets and the sizes of ligands are varying. They can be as small as ions or large polymers. Figure 2.2 shows an example of a pocket identified by the protein pocket identification tool GHECOM, that overlap with the ligand dysiherbaine bound to human glutamate receptor, GluR5 (PDB ID: 3FV1).

Generally, clinically approved drugs are classified into two broad classes (i) small molecules which typically comprise of <100 atoms with a molecular mass of <1,000 Da and (ii) biologics such as antibodies, modified nucleic acids, peptides, and vaccines [46]. The small

molecules usually traverse cellular membranes before they can reach intracellular target molecules and work on a relatively limited surface area. On the other hand, biologics generally interact with a large surface area that contains multiple interaction sites. They are also larger in sizes which restrict their mode of delivery and pose more challenges in drug design.



Figure 2.2 Pocket (orange sphere) identified by the protein pocket identification tool GHECOM that overlaps with the ligand dysiherbaine (red sticks) bound to human glutamate receptor, GluR5 (PDB ID: 3FV1). The figure was generated using PyMOL Molecular Graphics System [47].

2.6 Computational Methods for Identification of Protein Pockets

As the experimental identification of a binding site is not always feasible, as alternative computational protein pocket identification methods are required. Nowadays, many computational approaches for protein pocket identification have been developed which greatly accelerate drug discovery and protein designs. These protein pocket identification tools and algorithms fall into five different categories according to the methods applied (i) geometry based, (ii) energy based, (iii) evolution based, (iv) blind docking, and (v) combined approaches [48]. Geometry based methods can be grouped into three subcategories (i) grid system scanning, (ii) probe sphere filling, and (iii) based on the alpha shape theory [49,50]. In the following sections these methods are introduced in more details.

2.6.1 Geometry based Methods

Geometry based methods aim to identify solvent accessible regions that are located in surface cavities and clefts by the analysis of geometric criteria [48]. These approaches have been popular for years because they appear to exhibit good performance in the identification protein pockets. They were reported to predict almost 95% of the known binding sites. At a fast computational speed and they are robust in handling cases of structural variations or missing atoms/residues in the protein complexes [51]. Generally, geometry based methods are based

on the hypothesis that the ligand binding site are represented by the largest pocket [52–54]. Although these methods only considered the largest pocket, in reality this is not always the cases. Therefore, these methods have been further improved and new algorithms of prediction have been developed. In the following subsections the subcategories of geometry based methods are introduced in more detail.

2.6.1.1 Grid System Scanning

Grid based methods are a subcategory of geometry based methods that use a grid to define the molecular surface. These methods typically focus on the buriedness of grid points and the protein surface [54,55]. For example, DogSiteScorer is a grid based method. When tested on a dataset of 1069 structures, it achieved prediction accuracies more than 80% [56]. Generally, the initial step in the DoGSiteScorer process is the prediction of pockets on the protein surface based on the coordinates of the protein heavy atoms. A grid is spanned around the protein and grid points are labelled according to their spatial overlap with protein atoms. Then, the difference of Gaussian (DoG) filter is applied to the grid to identify the position of a cavity on the protein surface. Next, these positions are clustered to potential subpockets based on a density threshold. Finally, pockets are identified as the collection of merged neighbouring subpockets [48,56]. On the other hand, DEPTH [57,58], GHECOM [49], PocketDepth [59], and PocketPicker [55] also implement grid system scanning methods.

2.6.1.2 Probe Sphere Filling Methods

Probe sphere filling methods identify protein pockets by generating a set of probe spheres to fill cavities in a protein. Those regions containing the largest number of spheres are defined as pockets. In addition, these approaches use different types of probes such as (i) gap sphere [60], (ii) rotating probe, (iii) multiscale probe [49], (iv) the combination of big and small probes [61], and (v) probes placed tangential to triplets of protein atoms [62]. For example, SURFNET [60] and IsoMif [63] are representatives of these methods.

2.6.1.3 Alpha Shape Methods

In the1990s, the studies by Edelsbrunner and co-workers led to the development of new methods based on the alpha-shape theory [64,65]. The accuracy of the alpha-shape methods is influenced by the alpha values. The Automatic Protein Pocket Search (APROPOS) [66], the Computed Atlas of Surface Topology (CAST) [67], and Fpocket [68] are tools that implement alpha shape methods. For instance, CAST applied alpha shape methods to compute a triangulation of the protein surface atoms. This tool uses the discrete flow concept by allowing the small triangle flow to the neighbouring bigger triangles which act as "sinks" that later collect excess flow from neighbouring empty triangles. The collection of empty triangles is defined as the protein pocket [67].

2.6.2 Energy based Methods

Energy based methods rely on the assumption that the potential binding sites are characterized by binding energies which are different from the rest of the protein surface. Usually, these methods uses simple van der Waals (vdW) probes to locate the grid points around the protein surfaces and calculate interaction energies between the probe and a protein. The protein pockets are defined only by the energetic features. For examples, ProACT2 [69,70] and Q-SiteFinder [71] are tools that apply energy based methods to predict protein pockets.

2.6.3 Evolution based Methods

Evolution based methods also known as sequence based evolutionary conservation methods, are based on the assumption that functionally important residues of proteins are typically conserved during the evolution because of natural selection and these functional areas are mainly the protein pockets. Generally, the degree of conservation at each amino acid site is divided into two categories (i) slowly evolving sites which refer to evolutionarily conserved and (ii) rapidly evolving sites [72]. Based on this, several tools have been developed such as ConSeq [73], ConSurf [72], and a new version of ConSurf called Consurf 2010 [74] to identify the functional areas of unknown proteins by comparing their amino acid sequences to the already known amino acid sequences of proteins. These methods are fast, robust, and require only a protein sequence to predict the binding sites. Exceptions to this are the methods by de Rinaldis et al. 1998 [75] and siteFINDER|3D [76].

2.6.4 Blind Docking and Molecular Dynamic Methods

In addition to methods mentioned above, molecular dynamics (MD) and blind docking algorithms are also considered as useful methods to identify protein pockets. These approaches are most practical in cases where the ligand that binds in the target protein is known, but the binding site is unknown and requires the 3D-structure of the protein. Generally, MD calculations use two different approaches which are based on two situations (i) mobility of water molecules and (ii) long scale MD simulations to find the correct binding site of an already known active ligand. For example, a representative for this methods is molecular-docking binding site finding (MolSite) [77]. On the other hand, blind docking methods work by scanning the entire surface of the protein target to identify binding sites and modes of peptide ligands [78]. AnchorDock [79] is a representative of blind docking methods.

2.6.5 Combined Approaches

So far, apparently none of the methods described above performed in a way which can set as a standard or that is widely accepted as benchmarking method to identify protein pockets accurately. Several of the existing methods are often unsuccessful in certain types of cases, in which the algorithms are not able to take into account correctly or adequately all properties of the target site and did not work for all the entries in the dataset [80,81]. Furthermore, improving the existing algorithms does not necessarily can produce better predictions results. Thus the idea of extending the methods in other ways which aim to reduce the weaknesses of other algorithm appears promising. Consequently, the combination of two or more methods seems to be a good solution to improve the identification of protein pockets.

The first attempt of combined approaches was started by a study of Del Sol Mesa and co-workers (2003) [82]. They combined three separate evolution based methods but did not combine different types of approaches. Later, a study by Huang B et al., 2009 [83] introduced the MetaPocket tool which uses and combines the four methods LIGSITE^{cs} [84], PASS [62], Q-SiteFinder [71], and SURFNET [60] to predict protein pocket binding sites. The method was tested on two different datasets of 48 unbound/bound structures and 210 bound structures. The results showed a success rate of 70 to 75%. Two years later, Zhang et al. and co-workers introduced MetaPocket2.0 [85]. By applying these combined approaches to the same datasets, the results increased by 5% compared to the previous study.

2.7 Transcription Factors, Promoters, Motifs, and Position Weight Matrix (PWM)2.7.1 Transcription Factors and Promoters

The information within a gene is expressed by the cellular processes of transcription and translation. Transcription factors are regulatory DNA-binding proteins which play major roles in the regulation of gene expression. Previous studies noted that the larger total number of transcription factors reflects the larger number of genes in the genome of an organism which results in a larger size of the genome [86,87]. Single transcription factor or transcription factor complexes bind to the promoter region of the coding sequence which influencing RNA transcription. In addition, transcription factors can regulate gene expression by either activating or repressing gene transcription.

On the other hand, promoters are defined as the genomic regions that are located 5'upstream of the transcription start site (TSS) of genes. Promoters are the important element of expression vectors as they control the attachment of RNA polymerase and required to recruit the transcription initiation complexes and initiate transcription. There is no accurate definition of promoter length. Promoter binding is different in bacteria compared to eukaryotes as in bacteria RNA polymerase only needs the association of protein sigma factor to bind to the promoter. In contrast, eukaryotes require several transcription factors for the binding of RNA polymerase II to the promoter. Usually the complex consists of specific TFs, general TFs, cofactors, and RNA polymerase II [88]. Transcription does not only depend on co-factors and TFs and their ability to work together, but also on the structure of the chromatin. Additional mechanisms that also control gene expression are RNA interference, methylation and acetylation. The EPDnew database is one of the several publicly available resources of species-specific databases of experimentally validated promoters [89].

2.7.2 Motifs and position weight matrix (PWM)

Transcription factor binding sites (TFBSs) are the specific recognition sites in the DNA sequence for a given transcription factor. These binding sites are often referred to as occurrences of the motif for the corresponding TF. Nowadays, there are several approaches for the genome-wide detection of TFBSs such as (i) computational approaches which are based on the DNA sequence-based analysis and (ii) experimental approaches, which are nowadays based on chromatin immunoprecipitation (ChIP) and DNaseI HS-based technologies.

A position weight matrix (PWM) is a common way of representing patterns in biological sequences. It consists of a stack of letters representing each nucleotide at each position and the height of each letter is proportional to its value in the PWM. In addition, the sequence logos are regularly used to visualize count or frequency matrices. Figure 2.2 shows an example formation of PWM and the resulting sequence motif for eight known genomic binding sites in three genes HEM13, ANB1 and ROX1 from *Saccharomyces cerevisiae*. PWMs and sequence logos are available in several databases such as the JASPAR [90], HOCOMOCO [91], and TRANSFAC [92]. In this thesis, we identified the putative TFBS in the promoter regions of human STIM and ORAI genes. The description of STIM and ORAI genes are discussed in the Introduction section (4.1) of Chapter 4.



Figure 2.3 Example of formation of position weight matrix (PWM) and sequence motif. (A) Shown are eight known genomic binding sites in three *Saccharomyces cerevisiae* genes (HEM13, ANB1 and ROX1). (B) Frequencies matrix of nucleotides at each position. (C) Sequence logo used to visualize count at each position. (D) Sequence logo to represent the frequencies scaled relative using the information content at each position. (E) Energy normalized logo using relative entropy to adjust the GC content in *S.cerevisiae*. The figures were taken from [93,94].

GC content, guanine-cytosine content.

2.8 Gene Expression, Differential Gene Expression, and Co-expression Analysis

2.8.1 Gene Expression

The high-throughput sequencing technologies (HTS) are the most common approaches used in genomic studies which later involve statistical analysis to measure quantitative differences between experiments. It is important to analyze RNA expression levels (also known as RNA-seq data) to detect which genes that are differentially expressed across a group of samples.

Nowadays, there have been active efforts due to the advent of sequencing technologies with reduced costs that produce detailed profiling of gene expression levels, which are important in life sciences fields and clinical use [95]. Furthermore, the RNA-seq technology has been used to study complex transcriptomes and has assisted identification of levels of transcripts and isoforms, translocation events, sequence variations (for example SNPs) in the transcribed regions and post-transcriptional base modification.

Hence, various statistical approaches and tools have been developed such as differential expression analysis, random effects, gene set enrichment, gene set testing, and co-expression analysis to analyze the large datasets of genome-wide gene expression experiments. In this thesis, the differential gene expression analysis, co-expression analysis and differential interactions networks was performed using a dataset of breast invasive carcinoma (BRCA) obtained from the online data portal, The Cancer Genome Research (TCGA) (see Chapter 4).

2.8.2 Differential Gene Expression Analysis

The rapid growth of high-throughput technologies and publicly available datasets of RNA expression levels has motivated many studies to develop statistical algorithms that implement various approaches for normalization and differential gene expression analysis such as edgeR [96], DESeq [97], DESeq2 [98], PoisssonSeq, baySeq, and linear models for microarray data (limma). Generally, differential gene expression analysis of RNA-Seq data consists of three steps such as normalization of counts, parameter estimation of the statistical model and statistical tests for different expression.

2.8.2.1 DESeq, DESeq2, and edgeR Packages

DESeq, DESeq2 the extendable version of successful DESeq method and digital gene expression in R (EdgeR) are widely used Bioconductor packages for differential expression analysis of RNA-Seq and ChIP-Seq count data. These packages are very popular among user with biological background because they are easy to use, well documented and perform the best in replicated experiments. Generally, these statistical methods are based on the negative binomial distribution. Moreover, these packages have similar steps to perform differential analysis for count data. However they differ in several aspects such as (i) their look and feel, (ii) default normalization for example, edgeR applies the trimmed mean of M values and DESeq applies a relative log expression approach, and (iii) the application of dispersion estimate [99]. Typically, edgeR uses moderates individual dispersion estimates toward a trended-by-mean estimation [96]. A study by Dillies et al. 2013 [100] showed that the

normalization methods of edgeR and DESeq2 packages outperforms other approaches either in the case of expressed RNA repertoires that vary across biological conditions or in the presence of highly expressed genes. On the other hand, for the quality control step, clustering and Principal Component Analysis (PCA) can be used to assess the structure of the dataset.

In differential gene expression analysis, the False Discovery Rate (FDR) was used as a multiple testing correction approach. This was introduced by Benjamini and Hochberg (BH) in 1995 [101]. Multiple testing corrections adjust the p-values obtained from a large number of hypotheses testing to correct the occurrence of false positives. The FDR is defined as the expected proportion of falsely rejected null hypotheses among all rejected null hypotheses.

$$FDR = E\left(\frac{\text{number of falsely rejected null hypotheses}}{\text{number of rejected null hypotheses}}\right)$$
(2.1)

For example, a threshold of 0.02 FDR indicates that only two false positives are expected in 100 predictions. DESeq2 computes a q-value (FDR adjusted p-value) for each gene and uses it as the threshold to identify differentially expressed genes.

2.8.3 Co-expression Analysis

Generally, gene co-expression analysis is used to detect gene pairs that are coordinated in their expression profiles and to explore network characteristics of complex traits. In addition, gene co-expression network differential analysis is designed to assist biologists in many applications such as discovering protein-protein interaction relationships, predicting new gene functions, pathways, and identifying disease biomarkers or genes. In this network nodes represent genes and edges link two genes to show to what degree this pair of genes is co-expressed across several samples. The edges are based on correlation coefficients between each gene pair, where higher correlation means a higher probability of existing co-functionality between them. Recently, several computational approaches have been developed for the co-expression analysis and networks such as the Average Specific Connection, Differential Coexpression profile (DCp) [102], Differential Coexpression enrichment (DCe) [102], Differential Correlation in Expression for meta-module Recovery (DICER) [103], DiffCoEx [104], Log Ratio of Connections [105], and the Weighted Gene Co-Expression Network Analysis (WGCNA) [106]. For example, the popular R-package WGCNA has been used to analyze high dimensional data, distinct biological states and complex diseases. This approach identifies candidate genes relevant to a particular process of interest, construction of networks, module identifications, calculation of topological properties and visualizations [106]. The WGCNA package does not only concentrate on the individual gene expression, but it also focuses on modules of the genes which provide the relationships between modules [107]. In this thesis, we applied the WGCNA package for our co-expression analysis and differential interaction networks analysis.

Chapter 3

Composition of Overlapping Protein-Protein and Protein-Ligand Interfaces

This chapter is based on the following publication:

Ruzianisra Mohamed, Jennifer Degac, and Volkhard Helms. Composition of Overlapping Protein-Protein and Protein-Ligand Interfaces. PloS One. 2015 Oct 30;10 (10):e0140965. doi: 10.1371/journal.pone.0140965 [108].

My contribution was to write the manuscript, designed the research project and analyzed the results together with the co-authors Jennifer Degac and Volkhard Helms. The calculations were performed by me and Jennifer Degac.

Abstract

Protein-protein interactions (PPIs) play a major role in many biological processes and they represent an important class of targets for therapeutic intervention. However, targeting PPIs is challenging because often no convenient natural substrates are available as starting point for small-molecule design. Here, we explored the characteristics of protein interfaces in five non-redundant datasets of 174 protein-protein (PP) complexes, and 161 protein-ligand (PL) complexes from the ABC database, 436 PP complexes, and 196 PL complexes from the PIBASE database and a dataset of 89 PL complexes from the Timbal database. In all cases, the small molecule ligands must bind at the respective PP interface. We observed similar amino acid frequencies in all three datasets. Remarkably, also the characteristics of PP contacts and overlapping PL contacts are highly similar.

Keywords: Atomic contact, protein interface, protein-ligand interaction, protein-protein interaction, and amino acid composition.

3.1 Introduction

Protein-protein interactions (PPIs) play major roles in many biological processes such as bioenergetics, immune response, signal transduction, structural organization, and apoptosis [10,109]. Recently, PPIs also became a promising new target for therapeutic intervention. Unlike established pharmaceutical efforts that are directed, for example, at enzymes, G-protein coupled receptors (GPCR), or ion-channels, PPIs are challenging subjects because there are usually no convenient natural substrates that can be exploited as starting points for smallmolecule design. Moreover, the lack of information about particular interface residues determining the affinities and specificities at such interfaces makes it quite hard to design compounds that are capable of interfering with PPIs. Hence, there is a strong need to characterize the properties of protein interfaces that may also bind small-molecule ligands and the underlying molecular principles of contacts they are involved in.

The Protein Data Bank (PDB) [24] is the primary resource for elucidating the diversity of atomic contacts in protein-protein (PP) and protein-ligand (PL) interactions. Many statistical analyses of molecular interactions have been done based on this resource [1, 4–6]. Furthermore, some secondary databases that are derived from the PDB have been created to assist the integrated research on PP and PL interactions. Examples for this are the Timbal database (http://mordred.bioc.cam.ac.uk/timbal) which stores data of small molecules modulating protein–protein complexes [33], the Mother of All Database (MOAD) which contains data on ligand-protein binding (http://bindingmoad.org) [8-9], the 2P2I database of structures of PP complexes with known small molecule inhibitors (http://2p2idb.enrs-mrs.fr) [36], the Analysing Biomolecular Contacts (ABC) database (http://service.bioinformatik.uni-saarland.de/ABCSquareWeb/) [31], and the database of structurally defined protein interfaces named PIBASE (http://pibase.janelia.org/pibase2010/queries.html) [32]. One important aim in interface analysis is to identify properties which may distinguish binding residues from the rest of the protein surfaces.

Although protein-protein interfaces are rather large, planar and well packed depending [10,11], some parts of these interfaces termed overlap or bifunctional regions may bind both to small-molecule ligands and to proteins. The remaining regions of the interface which bind only to either protein or ligand are called non-overlap or monofunctional regions. Davis and Sali [113] found that bifunctional regions were enriched in tyrosine and tryptophan residues and depleted from alanine, isoleucine, leucine and valine when compared to monofunctional positions. Walter et al. [114] found for a different dataset that the overlap regions were mostly

found in pockets and some of their surfaces were exposed to the solvent. Koes and Camacho [115] used Small Molecular Inhibitor Starting Points (SMISPs) from PL and PP complexes in the PDB to train statistical classifiers for predicting such SMISPs.

In this study, we analyzed the residue-residue and atomic contact frequencies and propensities of five non-redundant datasets i) 174 protein-protein and ii) 161 protein-ligand complexes from Walter [114], iii) 436 protein-protein and iv) 196 protein-ligand complexes from the PIBASE database [32], and v) a dataset of 89 protein-ligand complexes from the Timbal database [33]. Our main research question was to find out whether small molecule ligands have similar physio-chemical features as protein binding interfaces when they bind at overlapping PP/PL binding interfaces and this was indeed found to be the case.

3.2 Material and Methods

3.2.1 Datasets

Non-redundant datasets from three different databases were used to investigate the composition of protein interfaces. The first pair of datasets consists of 174 PP complexes and 161 PL complexes compiled by Walter et al. [114] from the ABC database [31] (see Tables A and B in Supplementary Information Table 3.1). 25 entries of this PL dataset had been updated in the PDB in the meantime. We changed 22 previous ligand names to the current ligand names in the PDB files and removed 14 PDB files because they contain modified residues that were wrongly recognized as ligands before [114]. As described by Walter et al. [114], these complementary PP and PL datasets fulfill the following criteria: (i) PP: PL pairs represent pairs of complexes, where one protein may bind either a second protein or a small molecule ligand at the same interface, (ii) every pair of the dataset is represented as (Pi1, Pi2): (Pi3, Lj), where Pi1, Pi2 and Pi3 are three proteins and Lj is a small molecule ligand, (iii) Pi1 and Pi3 share at least 40% sequence identity, and (iv) the aligned positions in the binding interfaces of Pi1–Pi2 and Pi3 – Lj have at least two residues in common.

The same criteria of (Pi1, Pi2):(Pi3, Lj) pairs of PP and PL complexes from Walter et al., were then applied to the datasets of PP and PL complexes from the PIBASE database [32]. To avoid redundancy among these complexes, we clustered the PL complexes using the CD-Hit program [116,117] with the same sequence identity cut-off of 40%. Within a cluster, we selected the representative PP:PL pair with the highest identity score of the interface residues. Additionally, we discarded clusters which contained only sequences with fewer than 40 amino acids. The final pair of datasets comprises 436 PP complexes (Table C in Supplementary

Information Table 3.1) and 196 PL complexes (Table D in Supplementary Information Table 3.1).

Interactions where both interacting chains have > 90% sequence identity are defined as homodimer complexes and the remainder as heterodimer complexes. As a result, the PP complexes from the ABC dataset comprised 94 homodimer complexes and 80 heterodimer complexes (see Tables A and B in Supplementary Information Table 3.2). The PP complexes from the PIBASE dataset were grouped into 335 homodimer complexes and 101 heterodimer complexes (see Tables A and B in Supplementary Information Table 3.3).

The fifth dataset was extracted from the table of PDB entries in the Timbal database (see Table E in Supplementary Information Table 3.1). First, the 1695 entries in the current version of the Timbal database were filtered by removing complexes containing ligands that are annotated to act as stabilizers. Then, the CD-Hit program was applied to remove redundancy among the protein chains of the complexes with the sequence identity cut-off of 40%. We also eliminated clusters of proteins with fewer than 40 amino acids. This gave a final dataset of 89 protein-small molecule complexes.

Data from the ABC, PIBASE, and Timbal databases was retrieved by using MySQL queries, Java, Biojava [118] and analyzed with the R software (<u>http://www.R-project.org</u>).

3.2.2 Surface and Interface Residues

The solvent accessible surface area (SASA) was calculated using the NACCESS program [119]. As surface residues we considered those residues with a SASA value larger than zero. Labeled as interface residues were those residues that are within a radius of either 3 Å, 4 Å or 5 Å of any residue of the binding partner. Figure 3.1 shows a schematic diagram how we determined the interface and the remaining surface of PL complexes.



Figure 3.1 Schematic illustration of a PL complex illustrating the interface (black border) and the remaining surface regions. PL, protein-ligand.

3.2.3 Classification of the Amino Acids

The standard classification according to the Eisenberg hydrophobicity scale [120] was used to classify amino acids into four categories: hydrophobic (Ala, Ile, Leu, Met, Phe, Pro, Val), charged (Arg, Asp, Glu, Lys), polar (Cys, Asn, Gln, His, Ser, Thr, Trp, Tyr), and Gly.

3.2.4 Interface Residue Propensities

Residue interface propensities were calculated for the homodimeric and heterodimeric proteinprotein complexes of the ABC and PIBASE datasets and for the protein-ligand complexes of the ABC, PIBASE and Timbal datasets. These propensities give a measure of the relative importance of different amino acid residues in the interface, compared with the surface as a whole. The propensities were calculated with the following formula:

Interface residue propensity
$$AAj = \left(\frac{\sum \text{ interface residues of type } j}{\sum \text{ all interface residues}}\right) / \left(\frac{\sum \text{ surface residues of type } j}{\sum \text{ all surface residues}}\right)$$
(3.1)

An interface residue propensity of >1.0 indicates that a residue type occurs more frequently in interfaces than on the protein surface in general.

3.2.5 Contacts between Amino Acids of the Two Proteins

For every PP complex, we counted the observed number of contacts between amino acids of the first protein and amino acids of the second protein. A contact exists between two residues of these proteins if any residue of the first protein is within a distance threshold of 5.0 Å from the other protein. This was represented in a 20 x 20 table. From the 400 observed counts of amino acid pairs in the two datasets of protein-protein complexes, we derived normalized pair frequencies with the following formula:

$$Normalization = \frac{\left(\frac{\sum contacts of residue pair XY}{\sum all residue contacts}\right)}{\left(\frac{\sum observed X on surface}{\sum all surface residues of first protein}\right)\left(\frac{\sum observed Y on surface}{\sum all surface residues of second protein}\right)}$$
(3.2)

Here, XY is the number of observed contact pairs between residues X and Y across the interface, X is the count of amino acid X in the first protein and Y is the count of amino acid Y in the second protein.
3.2.6 Atom Contacts in Protein-protein and Protein-ligand Complexes

In protein-protein and protein-ligand complexes, we considered two surface atoms belonging to separate molecules to be in contact and labeled them as interface atoms if the distance between them is less than 5.0 Å. We counted contacts between all pairs of carbon (C), fluorine (F), nitrogen (N), oxygen (O), phosphorus (P), and sulfur (S) atoms resulting in 36 contact pairs. Then, the absolute counts were normalized as follows:

$$Normalization = \frac{\left(\frac{\sum contacts of atom pair AB}{\sum all atom contacts}\right)}{\left(\frac{\sum observed A on surface}{\sum all surface atom of ligand or protein}\right)\left(\frac{\sum observed B on surface}{\sum all surface atom of protein}\right)}$$
(3.3)

where A is the count of atom type A in the first protein (PP complexes) or protein (PL complexes), B is the count of atom type B in the second protein (PP complexes) or ligand (PL complexes) and, AB is the number of observed contact pairs between atom types A and B across the interface.

According to Higueruelo et al. [121], atom type contacts were grouped into polar and apolar contacts as follows: For protein-protein complexes, apolar contacts exist between C...C, C...S and S...S (not in Cys-Cys bridges). Polar contacts involve the pairs N...O, O...O, N...N, O...S and N...S (from Cys). For protein-ligand complexes, apolar contacts are C...C, and C...S pairs whereas polar contacts are formed by the pairs N...O, O...O, N...N, O...S, N...S, N...F, O...F, and S...F (from Cys).

3.2.7 Calculation of Polarity Ratio and Interface Atom Ratio

The polarity ratio (PR) is a simple measure of the polarity of the interface [122]. It was defined as the ratio of the number of polar atoms N, O, S at the interface to the sum of all C, N, O, S at the interface.

The interface atom ratio (IR) is a measure for the fraction of surface atoms that are located at the interface. It was calculated for the interfaces of protein-protein and protein-ligand complexes. Only the six heavy atom types C, N, O, S, P and F were considered in the calculation. IR is the ratio of the sum of all atoms at the interface to the sum of all atoms at the surface.

3.3 Results and Discussion

PPI interfaces are known to possess particular geometric and physicochemical characteristics, see e.g. [10,123–125]. Comparing these features of protein interfaces to those of overlapping protein-ligand interfaces should aid in targeting protein-protein interaction sites. Here, we used the ABC, PIBASE and Timbal databases as data sources for protein interfaces and surfaces. All three databases are secondary database that are derived from the PDB. However, due to the different way of identifying overlapping PP/PL pairs, the direct overlap between the three nonredundant datasets derived from them is fairly small. We believe that this may have resulted from the clustering with the CD-Hit program that selected different cluster representatives in each case. We found only the following redundant PP complexes 1AB8 (B-A), 1AZZ (A-C), 1BMF (F-B), 1EYS (H-M), 1RQ8 (A-E), 1SGF (G-B) from the ABC dataset and 1AB8 (A-B), 1AZZ (C-A), 1BMF (C-D), 1EYS (M-C), 1RQ8 (E-A), 1SGF (G-Z) from the PIBASE dataset. Furthermore, both datasets share the following lists of PDB IDs with same chain interactions 1DPJ (A-B), 1POS (H-E), and 2G2U (A-B). Similarly, there are few redundancies between datasets of PL complexes from ABC and PIBASE, namely 1C50 (A-CHI), 1KYN (A-KTP), 1LBC (A-CYZ) and 1M2Z (A-BOG), respectively. There is also one overlapping member between the datasets of PL complexes from PIBASE and Timbal, namely the PDB ID 1AB8 (A-FOK). Figure 3.2 summarizes the workflow of the analysis of the five datasets. The fraction of homodimers and heterodimers in the datasets derived from ABC and PIBASE are 54%: 46% and 77%: 23%, respectively.





PP, protein-protein; PL, protein-ligand; PR, polarity ratio; IR, interface atom ratio.

3.3.1 Amino Acid Composition and Protein Interfaces Propensity

Figures 3.3 and 3.4 show the percentage frequencies and propensities of amino acids at the interfaces of homodimeric and heterodimeric PP complexes from the ABC and PIBASE datasets, respectively. Figure 3.5 shows the percentage frequencies and propensities of amino acids at the protein interfaces of the PL complexes from the ABC, PIBASE and Timbal datasets, respectively. Previous studies showed that protein-protein interfaces have unique characteristics that distinguish them from non-interface portions of protein surfaces [123,126,127].

By grouping the amino acids according to the Eisenberg hydrophobicity scale (see methods) we found that, hydrophobic amino acids account for 38.06% (ABC-P1-homo), 38.87% (ABC-P2-homo), 38.81% (PIB-P1-homo) and 38.75% (PIB-P2-homo) at interfaces of homodimeric PP complexes compared to 35.60% (ABC-P1-hetero), 36.11% (ABC-P2-hetero), 37.94% (PIB-P1-hetero) and 36.38% (PIB-P2-hetero) at interfaces of heterodimeric PP complexes (Figures 3.3A and 3.4A). This matches the general finding e.g. of Jones and Thornton who stated that homodimer complexes are more hydrophobic [10].

At interfaces of both homodimeric and heterodimeric PP complexes from the ABC and PIBASE datasets, alanine, valine, and lysine residues are underrepresented with propensities lower than 1.0 (Figures 3B and 4B). One hydrophobic amino acid (leucine), one charged amino acid (lysine) and two polar amino acids (glutamine and threonine) have higher propensities at interfaces of homodimer complexes than at interfaces of heterodimer complexes of the ABC dataset. In the PIBASE dataset, four hydrophobic amino acids (alanine, leucine, proline and valine), one polar amino acid (threonine) and glycine have higher propensities in homodimer complexes than in heterodimer complexes.



Figure 3.3 Percentage frequencies and propensities of amino acid residues at protein interfaces of PP complexes from the ABC dataset. (A) Percentage frequencies of amino acid residues at protein interfaces. (B) Propensities of amino acid residues at protein interfaces.

PP, protein-protein; ABC, ABC dataset; homo, homodimeric PP interface; hetero, heterodimeric PP interfaces; P1, protein interface of the first protein (Pi1); P2, protein interface of the second protein (Pi2).



Figure 3.4 Percentage frequencies and propensities of amino acid residues at protein interfaces of PP complexes from the PIBASE dataset. (A) Percentage frequencies of amino acid residues at protein interfaces. (B) Propensities of amino acid residues at protein interfaces. PP, protein-protein; PIB, PIBASE dataset; homo, homodimeric PP interface; hetero, heterodimeric PP interfaces; P1, protein interface of the first protein (Pi1); P2, protein interface of the second protein (Pi2).

As expected, hydrophobic and polar residues make up the largest portion of protein interfaces. In fact, this is one of the challenges for targeting PPIs with small molecules as the contact surfaces between proteins typically involve many hydrophobic and polar interactions distributed over a large interface with buried area of ~1500 – 3000 Å² [128]. According to the classification by Eisenberg, the fractions of hydrophobic, polar, charged and glycine residues are 36.95%, 33.38%, 22.11%, 7.56% for the first protein (Pi1), 37.70%, 32.48%, 22.35%, 7.46% for the second protein (Pi2) of the PP complexes from the ABC dataset, 38.60%, 30.93%, 24.09%, 6.38% for the first protein (Pi1), and 38.20%, 31.35%, 24.22%, 6.23% for the second protein (Pi2) of the PP complexes from the PIBASE dataset. Although there are minor differences between the two datasets (slightly more charged and fewer glycine residues in the PIBASE dataset), we found the composition to be overall remarkably similar.

At the interfaces of both homodimeric and heterodimeric PP complexes from the ABC and PIBASE datasets, the frequencies of methionine and tryptophan at protein interfaces are at most 3.07%. However, both amino acids have normalized interface propensities clearly larger than one, suggesting that these residues play important roles and thus occur more frequently at protein interfaces rather than elsewhere on the protein surface. Overall, tryptophan, tyrosine and arginine each have propensities above 1.0 at both protein interfaces of homodimeric and heterodimeric PP complexes from the ABC and PIBASE datasets. This reflects that aromatic amino acids and arginine play important roles in protein interfaces, which is a well-known fact. For example, Bogan and Thorn [129] reported that hotspot regions at protein interfaces are enriched in tryptophan, tyrosine and arginine. Also, Jones, Marin and Thornton [130] found that hydrophobic residues including tryptophan and tyrosine as well as arginine are moderately enriched at protein interfaces compared to the whole surface. Jones and Thornton [10] reported that with the exception of methionine, all hydrophobic residues show a greater preference for the interfaces of homodimers than for those of heterocomplexes. Based on our analysis, only leucine is clearly enriched at homodimer interfaces. Janin, Bahadur and Chakrabarti [125] wrote that relative to the accessible protein surface, the interfaces are depleted in glutamic acid, aspartic acid, and lysine, and enriched in methionine, tyrosine and tryptophan. Our findings are in good agreement with this. In our case, the enriched category also includes phenylalanine, histidine and arginine. The underrepresented category also includes alanine, proline and valine. Talavera et al. [131] provided a rather recent compilation of amino acid frequencies and propensities, separately for homomeric and heterodimeric PP complexes. A possible concern about their work is that they applied a rather generous homology threshold of 80% identity. They found tyrosine, tryptophan, methionine, cysteine, phenylalanine, leucine, valine and isoleucine to be enriched at the interfaces of homo-complexes. In the case of hetero-complexes, cysteine fell out from this list. On the other hand, lysine, asparagine, aspartic acid and glutamic acid were underrepresented in homo-complexes. The same ones plus serine and glycine were found for hetero-complexes.

The distributions of the percentage frequencies and propensities of amino acids at the protein interfaces of the PP datasets derived from ABC and PIBASE were compared with the non-parametrical Friedman test as the datasets do not have a normal distribution. As suggested by the graphical representation of Figures 3.3 and 3.4, the ABC and PIBASE datasets do not differ significantly (percentage frequencies, p-value = 0.99 and propensities, p-value = 0.97).

The fractions of hydrophobic, polar, charged and glycine residues at protein binding interfaces of PL complexes are 34.08%, 36.97%, 20.31%, 8.64% (ABC dataset), 38.25%, 35.39%, 18.12%, 8.24% (PIBASE dataset) and 42.32%, 32.61%, 18.60%, 6.47% (Timbal dataset), see Figure 3.5A. Compared to PP interfaces, the ligand-contacting protein interfaces of the Timbal dataset contain about 5% more hydrophobic residues, and about 5% fewer charged residues. In contrast, the ligand-contacting protein interfaces from the ABC and PIBASE datasets contain 3-4% more polar residues than PP interfaces and 3-4% less charged residues.

In the PL complexes of the ABC dataset, the five amino acids with the highest propensities found at protein interfaces are cysteine (2.20), tryptophan (2.18), histidine (1.75), tyrosine (1.74), and phenylalanine (1.47). In the PL complexes of the PIBASE dataset, the most enriched ones are tryptophan (2.25), tyrosine (1.93), phenylalanine (1.92), histidine (1.89), and methionine (1.66). In the PL complexes of the Timbal dataset, methione has the highest propensity of 1.85, followed by phenylalanine (1.78), tryptophan (1.78), histidine (1.54) and tyrosine (1.53), respectively. In all datasets of PL complexes, tryptophan, phenylalanine, histidine, and tyrosine are found most often at the protein interfaces (Figure 3.5B) complemented by either cysteine (ABC) or methionine (PIBASE, Timbal).

The distributions of percentage frequencies and propensities of amino acids acids at the protein interfaces in the datasets derived from ABC, PIBASE and Timbal did not differ significantly (percentage frequencies, p-value = 0.86 and propensities, p-value = 0.96, Friedman rank sum test).



Figure 3.5 Percentage frequencies and propensities of amino acids residues at protein interfaces of PL complexes from the ABC, PIBASE and Timbal datasets. (A) Percentage frequencies of amino acids residues at protein interfaces. (B) Propensities of amino acids residues at protein interfaces. PL, protein-ligand; PL-ABC, PL complexes from the ABC dataset; PL-PIBASE, PL complexes from the PIBASE dataset; PL-Timbal, PL complexes from the Timbal dataset.

3.3.2 Amino Acid Contacts

The propensities of amino acid contacts in PP complexes between amino acids of the first protein (Pi1) and amino acids of the second protein (Pi2) were obtained by counting the absolute number of contacts and normalizing this number against the appearance probability of the two involved residues at the surface. In Figures 3.6 and 3.7, the propensity values were log2 transformed to ensure a balanced view of over- and under-representation. Contacts with high propensities were observed among residues pairs of different polarity types. In PP complexes from the ABC dataset, the five most over-represented interactions were found between the pairs of tryptophan (6.32), cysteine (4.66), phenylalanine (3.61) and histidine (3.50) as well as between tryptophan and phenylalanine (3.36), see Figure 3.6. In PP complexes from the PIBASE dataset, the five most over-represented interactions were pairs of tryptophan (7.50), methionine (4.34), phenylalanine (3.96), tyrosine (3.57), and cysteine (3.43), see Figure 3.7. These results are consistent with previous studies of protein-protein interfaces that reported

an enrichment of contacts between cysteine, hydrophobic contacts and aromatic contacts [123, 126, 132-134]. Further studies noticed that besides disulfide bonds and hydrophobic interactions, also salt-bridges contribute to stabilizing protein-protein interactions [126,132–134]. In our analysis, contacts between lysine and negatively charged amino acids (Asp, Glu) are only mildly enriched (propensity 1.23 on average), whereas those between arginine and either Asp or Glu are about two-fold enriched (2.06), see Tables A and B in S5 File, what reflects the enriched of arginine at protein interfaces. The propensities of amino acid contacts between amino acids of the first protein (Pi1) and amino acids of the second protein (Pi2) in PP complexes between datasets from the ABC and the PIBASE did not differ statistically significantly (p-value = 0.76, Wilcoxon signed rank test).



Figure 3.6 Amino acid pairing propensities (in log₂-format) for interfaces of PP complexes from the ABC dataset. PP, protein-protein.



Figure 3.7 Amino acid pairing propensities (in log₂-format) for interfaces of PP complexes from the PIBASE dataset. PP, protein-protein.

Based on the counts of amino acids, we computed the average number of amino acid residues at the interfaces of the two proteins Pi1 and Pi2 of PP complexes and the Pi3 protein of PL complexes using three different atom distances (3Å, 4Å and 5Å). At the distance threshold of 3Å, both interfaces at Pi1 and Pi2 contain less than 10 amino acids on average. For thresholds of 4Å and 5Å, the average size of the protein interfaces is 26.22 (ABC dataset) and 38.69 amino acids (PIBASE dataset) (Table 3.1).

Table 3.1 The average number with standard deviation of amino acid residues at the interfaces of PP complexes in the ABC and PIBASE datasets.

| | PP complexes | | | | | | | | |
|---------------|---------------|--------------|----------------|---------------|--|--|--|--|--|
| | ABC da | ataset | PIBASE dataset | | | | | | |
| Atom distance | Pi1 Pi2 | | Pi1 | Pi2 | | | | | |
| 3 Å | 7.67 ± 6.68 | 7.49 ± 6.85 | 9.61 ± 15.62 | 9.53 ± 15.52 | | | | | |
| 4 Å | 27.17 ± 19.18 | 26.22 ±19.7 | 31.13 ± 24.53 | 30.76 ± 24.61 | | | | | |
| 5 Å | 34.52 ± 23.31 | 32.8 ± 24.11 | 38.69 ± 27.97 | 38.09 ± 27.94 | | | | | |

PP, protein-protein; Pi1, first protein; Pi2, second protein.

Table 3.2 shows the average number of residues at the interfaces of Pi3 in PL complexes from the ABC, PIBASE and Timbal datasets. At the distance threshold of 3 Å, the average size of the interfaces is less than 3 amino acids for all datasets. At 4 Å and 5 Å atom distances, the average sizes of the interfaces are between 6.31 amino acids (ABC dataset) and 13.54 amino acids (Timbal dataset). Although the PL interfaces from the ABC dataset are clearly smaller than those from the PIBASE and Timbal datasets, the average ligand size in the ABC dataset (20.48 atoms without hydrogen atoms) is only slightly smaller than the average ligand size in the Timbal dataset (21.53 atoms) and in the PIBASE dataset (21.42 atoms), respectively.

Table 3.2 The average number with standard deviation of amino acid residues at the interfaces of PL complexes in the ABC, PIBASE and Timbal datasets.

| | PL complexes | | | | | | | |
|---------------|--------------|----------------|----------------|--|--|--|--|--|
| Atom distance | ABC dataset | PIBASE dataset | Timbal dataset | | | | | |
| 3 Å | 1.64 ± 1.93 | 2.58 ± 2.08 | 2.54 ± 2.52 | | | | | |
| 4 Å | 6.31 ± 4.66 | 10.04 ± 4.39 | 9.99 ± 6.32 | | | | | |
| 5 Å | 8.84 ± 5.79 | 13.43 ± 5.63 | 13.54 ± 8.06 | | | | | |

PL, protein-ligand.

3.3.3 Atomic Contacts in Protein-protein and Protein-ligand Complexes

In this section, we analyzed the atomic contacts in the datasets of PP and PL complexes. For atom pairs between the first and second proteins (Pi1–Pi2) in PP complexes and between protein and ligand (Pi3 – Lj) in PL complexes, we counted contacts of less or equal to 5 Å between six types of heavy atoms, namely carbon (C), flourine (F), nitrogen (N), oxygen (O), phosphorus (P) and sulfur (S). This resulted in 36 atomic pair contacts. Table 3.3 lists the appearance frequency of these 36 atomic contact types in PP and PL complexes from the ABC, PIBASE and Timbal datasets. In all datasets, the most frequent contacts are C...C (> 41%), O...C (> 10%), C...O (> 8%), and C...N (>7%).

Chen and Kurgan [135] previously characterized the binding interfaces of proteins with small molecules, irrespective of whether they also bind to other proteins. As expected, interactions with organic molecules are dominated by van der Waals contacts, hydrogen bonds, and covalent contacts, whereas those with charged species also involve electrostatic interactions. Hakulinen et al. [136] argued that small molecules frequently contact phenylalinine, histidine, tyrosine and tryptophan residues of proteins because their aromatic ring carbons prefer other aromatic carbons. Both findings match well with the results of this analysis. The atomic contacts in PP complexes of the ABC and PIBASE datasets did not differ

significantly (p-value = 0.76, Wilcoxon signed rank test). Also the frequencies of the atomic contacts between the PL complexes of the ABC, PIBASE and Timbal datasets did not differ significantly (p-value = 0.11, Friedman rank sum test).

Tables 3.4 and 3.5 list the percentage frequencies and normalized propensities of apolar, polar and other atomic contacts in PP complexes and PL complexes, respectively. The content of apolar contacts (45.52% for the PP complexes in the ABC dataset and 45.25% for the PIBASE dataset) and of polar contacts (13.85% vs 13.70%) is highly similar between the two PP datasets. In contrast, the PL complexes of the PIBASE dataset (46.45%) contained more apolar contacts than the Timbal dataset (44.84%) and the ABC dataset (43.04%). Concerning polar contacts in PL complexes, the Timbal dataset (14.71%) and the ABC dataset (14.48%) contain more such contacts than the PIBASE dataset (12.95%). Overall, the differences of the normalized propensities seem minor, among the PP and PL datasets, as well as between PP and PL datasets, which agrees with the findings of [135]. In all datasets, C-C contacts are slightly overrepresented (1.04 to 1.11 times the randomly expected number of contacts). N-N contacts are always more frequent (1.07 to 1.32) than O-O contacts (0.70 to 0.95).

| | | PP comple | xes | PL complexes | | | |
|--------------------|--------------------|-----------|---------|--------------|---------|---------|--|
| | | ABC | PIBASE | ABC | PIBASE | Timbal | |
| [] | | dataset | dataset | dataset | dataset | dataset | |
| Atom1 ^a | Atom2 ^b | % | % | % | % | % | |
| С | С | 44.08 | 44.20 | 41.24 | 44.57 | 43.49 | |
| С | Ν | 10.82 | 10.95 | 8.50 | 9.03 | 7.05 | |
| С | 0 | 10.42 | 10.68 | 12.63 | 11.85 | 8.88 | |
| С | S | 0.85 | 0.57 | 0.90 | 0.88 | 1.00 | |
| С | Р | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | |
| С | F | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | |
| Ν | С | 8.29 | 8.12 | 5.81 | 5.16 | 5.51 | |
| Ν | Ν | 2.76 | 2.70 | 1.58 | 1.25 | 1.30 | |
| Ν | 0 | 2.85 | 2.83 | 2.18 | 1.77 | 1.57 | |
| Ν | S | 0.14 | 0.08 | 0.17 | 0.09 | 0.13 | |
| Ν | Р | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | |
| Ν | F | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | |
| 0 | С | 10.93 | 11.17 | 13.87 | 12.78 | 15.03 | |
| 0 | Ν | 4.48 | 4.40 | 4.89 | 4.38 | 6.58 | |
| 0 | 0 | 3.37 | 3.45 | 4.79 | 4.43 | 4.60 | |
| 0 | S | 0.17 | 0.14 | 0.18 | 0.25 | 0.22 | |
| 0 | Р | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | |
| 0 | F | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | |
| S | С | 0.56 | 0.47 | 0.90 | 1.00 | 0.35 | |
| S | Ν | 0.14 | 0.10 | 0.26 | 0.26 | 0.08 | |
| S | 0 | 0.11 | 0.12 | 0.31 | 0.23 | 0.04 | |
| S | S | 0.03 | 0.02 | 0.04 | 0.04 | 0.01 | |
| S | Р | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | |
| S | F | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | |
| Р | С | 0.00 | 0.00 | 0.59 | 0.45 | 1.64 | |
| Р | Ν | 0.00 | 0.00 | 0.50 | 0.28 | 1.18 | |
| Р | 0 | 0.00 | 0.00 | 0.19 | 0.22 | 0.40 | |
| Р | S | 0.00 | 0.00 | 0.02 | 0.01 | 0.00 | |
| Р | Р | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | |
| Р | F | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | |
| F | С | 0.00 | 0.00 | 0.32 | 0.77 | 0.75 | |
| F | Ν | 0.00 | 0.00 | 0.04 | 0.11 | 0.09 | |
| F | 0 | 0.00 | 0.00 | 0.08 | 0.15 | 0.10 | |
| F | S | 0.00 | 0.00 | 0.01 | 0.04 | 0.00 | |
| F | Р | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | |
| F | F | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | |
| Total | | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | |

Table 3.3 The percentage frequencies of the 36 atomic contact types in PP and PL complexes.

PP, protein-protein; PL, protein-ligand.

^aFor PP complexes, atom1 belongs to the first protein and for PL complexes, atom1 belongs to the protein.

^bFor PP complexes, atom2 belongs to the second protein and for PL complexes, atom2 belongs to the ligand.

| | | PP complex | kes |
|---------------------|----------------------------|-----------------|----------------|
| | | ABC dataset | PIBASE dataset |
| Apolar contacts: | CC | 44.08 (1.10) | 44.20 (1.10) |
| | CS | 0.85 (2.53) | 0.57 (1.91) |
| | SC | 0.56 (1.82) | 0.47 (1.63) |
| | SS (not in Cys-Cys bridge) | 0.03 (10.94) | 0.01 (6.60) |
| | Total | 45.52 | 45.25 |
| Polar contacts: | NO | 2.85 (0.91) | 2.83 (0.90) |
| | 0N | 4.48 (1.40) | 4.40 (1.40) |
| | 00 | 3.37 (0.70) | 3.45 (0.73) |
| | NN | 2.76 (1.31) | 2.70 (1.29) |
| | 0S | 0.16 (1.40) | 0.14 (1.38) |
| | SO | 0.11 (1.05) | 0.12 (1.24) |
| | N…S (from Cys) | 0.05 (0.68) | 0.02 (0.35) |
| | SN (from Cys) | 0.07 (0.93) | 0.03 (0.50) |
| | Total | 13.85 | 13.7 |
| Other contacts: | CN | 10.83 (1.19) | 10.95 (1.20) |
| | NC | 8.28 (0.91) | 8.12 (0.88) |
| | CO | 10.42 (0.76) | 10.68 (0.78) |
| | 0C | 10.93 (0.78) | 11.17 (0.81) |
| | N…S (S not from Cys) | 0.09 (1.14) | 0.06 (0.86) |
| | S…N (S not from Cys) | 0.08 (1.08) | 0.07 (1.05) |
| | SS (in Cys-Cys bridge) | 0.002 (0.84) | 0.003 (1.39) |
| | CP/PC | 0.00 (0.00) | 0.00 (0.00) |
| | CF/FC | 0.00 (0.00) | 0.00 (0.00) |
| | NP/PN | 0.00 (0.00) | 0.00 (0.00) |
| | NF/FN | 0.00 (0.00) | 0.00 (0.00) |
| | 0P/P0 | 0.00 (0.00) | 0.00 (0.00) |
| | 0F/FO | 0.00 (0.00) | 0.00 (0.00) |
| | SP/PS | 0.00 (0.00) | 0.00 (0.00) |
| | SF/FS | 0.00 (0.00) | 0.00 (0.00) |
| | PP | 0.00 (0.00) | 0.00 (0.00) |
| | PF/FP | 0.00 (0.00) | 0.00 (0.00) |
| | FF | 0.00 (0.00) | 0.00 (0.00) |
| | Total | 40.63 | 41.05 |
| | Grand Total | 100 | 100 |

Table 3.4 Percentage frequencies (with normalized propensity values in parentheses) of apolar, polar and other atomic contacts of PP complexes from the ABC and PIBASE datasets.

PP, protein-protein.

| | | PL complexes | | |
|------------------|----------------------|--------------|----------------|----------------|
| | | ABC dataset | PIBASE dataset | Timbal dataset |
| Apolar contacts: | CC | 41.24 (1.04) | 44.57 (1.04) | 43.49 (1.11) |
| | CS | 0.90 (2.86) | 0.88 (2.82) | 1.00 (2.60) |
| | SC | 0.90 (1.26) | 1.00 (1.14) | 0.35 (0.99) |
| | Total | 43.04 | 46.45 | 44.84 |
| Polar contacts: | NO | 2.18 (1.20) | 1.77 (1.02) | 1.57 (0.94) |
| | ON | 4.89 (1.31) | 4.38 (1.39) | 6.58 (1.68) |
| | 00 | 4.79 (0.85) | 4.43 (0.95) | 4.6 (0.80) |
| | NN | 1.58 (1.32) | 1.25 (1.07) | 1.30 (1.14) |
| | 0S | 0.18 (1.45) | 0.25 (2.50) | 0.22 (1.33) |
| | SO | 0.31 (1.22) | 0.23 (0.77) | 0.04 (0.36) |
| | NS | 0.17 (4.07) | 0.09 (1.24) | 0.13 (2.75) |
| | SN | 0.26 (1.59) | 0.26 (1.30) | 0.08 (0.97) |
| | NF | 0.00 (0.00) | 0.00 (0.00) | 0.00 (0.00) |
| | FN | 0.04 (0.56) | 0.11 (0.46) | 0.09 (0.88) |
| | 0F | 0.00 (0.00) | 0.00 (0.00) | 0.00 (0.00) |
| | FO | 0.08 (0.80) | 0.15 (0.43) | 0.10 (0.70) |
| | S…F (S from Cys) | 0.00 (0.00) | 0.00 (0.00) | 0.00 (0.00) |
| | FS (S from Cys) | 0.00 (0.00) | 0.02 (3.21) | 0.00 (0.00) |
| | Total | 14.48 | 12.95 | 14.71 |
| Other contacts: | CN | 8.50 (0.93) | 9.03 (0.91) | 7.05 (0.77) |
| | NC | 5.81 (1.12) | 5.16 (1.03) | 5.51 (1.13) |
| | CO | 12.63 (0.91) | 11.85 (0.81) | 8.88 (0.66) |
| | OC | 13.87 (0.86) | 12.78 (0.94) | 15.03 (0.90) |
| | СР | 0.00 (0.00) | 0.00 (0.00) | 0.00 (0.00) |
| | PC | 0.59 (0.57) | 0.45 (0.94) | 1.64 (1.07) |
| | CF | 0.00 (0.00) | 0.00 (0.00) | 0.00 (0.00) |
| | FC | 0.32 (1.14) | 0.77 (0.77) | 0.75 (1.77) |
| | NP | 0.00 (0.00) | 0.00 (0.00) | 0.00 (0.00) |
| | PN | 0.50 (2.08) | 0.28 (2.52) | 1.18 (3.30) |
| | 0P | 0.00 (0.00) | 0.00 (0.00) | 0.00 (0.00) |
| | PO | 0.19 (0.51) | 0.22 (1.33) | 0.40 (0.77) |
| | SS | 0.04 (6.30) | 0.04 (6.34) | 0.01 (2.09) |
| | SP | 0.00 (0.00) | 0.00 (0.00) | 0.00 (0.00) |
| | PS | 0.02 (2.18) | 0.01 (2.90) | 0.00 (0.00) |
| | S…F (S not from Cys) | 0.00 (0.00) | 0.00 (0.00) | 0.00 (0.00) |
| | FS (S not from Cys) | 0.01 (5.40) | 0.01 (1.84) | 0.00 (0.00) |
| | PP, PF/FP and FF | 0.00 (0.00) | 0.00 (0.00) | 0.00 (0.00) |
| | Total | 42.48 | 40.6 | 40.45 |
| | Grand Total | 100 | 100 | 100 |

Table 3.5 Percentage frequencies (with normalized propensity values in parentheses) of apolar, polar, and other atomic contacts of PL complexes from the ABC, PIBASE, and Timbal datasets.

3.3.4 Polarity Ratio and Interface Atom Ratio

Then, we analyzed the polarity ratio (PR), namely the fraction of polar N, O, S atoms at the interface areas of both PP and PL complexes. The interface atom ratio (IR) indicates the fraction of surface atoms that are involved in protein contacts at the interface. As mention before, the interface areas were defined as those residues that are closer than 3 Å (or 4 Å and 5 Å) to at least one residue from the binding partner. Both IR and PR were computed for the datasets of PP and PL complexes from the ABC, PIBASE, and Timbal datasets.

At 3 Å distance threshold, the differences in IR and PR ratios are not representative because only the shortest-distance contacts are considered. For example, when a 3 Å cut-off is used, most carbon atoms are not considered as part of the interfaces as this short distance is shorter than twice the van der Waals radius of carbon (1.7 Å) [137]. Table 3.6 shows that, as expected, only very small differences were observed when computing PR and IR of PP complexes between the first protein (Pi1) and the second protein (Pi2), as both of them exhibit similar characteristics at binding interfaces. For the larger cut-off distances (4 Å and 5 Å), the polarity ratio (PR) decreases quickly because now all carbon atoms at the surface are included. On the other hand, the interface atom ratio (IR) of 8.0% (4 Å) and 14.0% (5 Å) shows that, expectedly, only a small fraction of the protein surface atoms are included in the interface.

| Table 3.6 | Interface | atom | ratio | (IR) | and | polarity | ratio | (PR) | (with | standard | deviations | in |
|------------|--------------|-------|-------|------|--------|------------|-------|---------|-------|------------|------------|----|
| parenthese | s) for inter | faces | of PP | comp | olexes | s from the | e ABC | C and I | PIBAS | E datasets | 5. | |

| | | PP c | complexes from ABC dataset | n the | PP complexes from the PIBASE dataset | | | | |
|----|-----|--------------|-------------------------------|--------------|---|--------------|--------------|--|--|
| | | | Atom distance | 9 | Atom distance | | | | |
| | | 3 Å | 4 Å | 5 Å | 3 Å | 4 Å | 5 Å | | |
| IR | | | | | | | | | |
| | Pi1 | 0.01 (±0.01) | 0.08 (±0.05) | 0.14 (±0.08) | 0.01 (±0.07) | 0.08 (±0.08) | 0.13 (±0.09) | | |
| | Pi2 | 0.01 (±0.01) | 0.09 (±0.08) | 0.15 (±0.13) | 0.01 (±0.07) | 0.08 (±0.08) | 0.13 (±0.10) | | |
| PR | | | | | | | | | |
| | Pi1 | 0.87 (±0.22) | 0.38 (±0.07) | 0.34 (±0.06) | 0.72 (±0.20) | 0.37 (±0.06) | 0.34 (±0.05) | | |
| | Pi2 | 0.85 (±0.22) | 0.37 (±0.06) | 0.34 (±0.05) | 0.71 (±0.20) | 0.37 (±0.06) | 0.34 (±0.05) | | |

IR, interface atom ratio; PR, polarity ratio; PP, protein-protein; Pi1, first protein; Pi2, second protein.

Table 3.7 lists the IR and PR ratios of 161 PL complexes from the ABC dataset, 196 PL complexes from the PIBASE dataset, and 89 PL complexes from the Timbal dataset. At the distance threshold of 3 Å, almost no ligands atoms are considered as interfacial atoms whereas the opposite is the case for 5 Å where 93% (PIBASE) and 94% (Timbal) of the ligand atoms are considered as interfacial atoms compared to 78% for ABC. This is suggesting that the

PIBASE and Timbal ligands bind more flat on the protein surfaces and/or bind deeper into pockets on the protein surface than the ABC ligands. Finally, the polarity ratios of the proteins in the PL dataset are comparable to the proteins in the PP dataset.

| | | PL complexes from the ABC dataset | | | PL complexes from the PIBASE dataset | | | PL complexes from the Timbal dataset | | |
|----|-----|--------------------------------------|---------|---------|---|---------|---------|---|---------|---------|
| | | Atom distance | | | Atom distance | | | Atom distance | | |
| | | 3 Å | 4 Å | 5 Å | 3 Å | 4 Å | 5 Å | 3 Å | 4 Å | 5 Å |
| | | 0.002 | 0.01 | 0.03 | 0.003 | 0.02 | 0.03 | 0.003 | 0.02 | 0.03 |
| IR | Pi3 | (±0.004) | (±0.02) | (±0.03) | (±0.003) | (±0.01) | (±0.02) | (±0.003) | (±0.02) | (±0.03) |
| | | 0.11 | 0.55 | 0.78 | 0.13 | 0.74 | 0.93 | 0.16 | 0.75 | 0.94 |
| | Lj | (±0.14) | (±0.27) | (±0.25) | (±0.17) | (±0.25) | (±0.20) | (±0.14) | (±0.21) | (±0.19) |
| | | 0.83 | 0.38 | 0.35 | 0.85 | 0.38 | 0.34 | 0.86 | 0.36 | 0.32 |
| PR | Pi3 | (±0.45) | (±0.20) | (±0.14) | (±0.36) | (±0.14) | (±0.12) | (±0.38) | (±0.16) | (±0.12) |
| | | 0.76 | 0.38 | 0.35 | 0.79 | 0.33 | 0.31 | 0.84 | 0.38 | 0.36 |
| | Lj | (±0.45) | (±0.23) | (±0.18) | (±0.38) | (±0.17) | (±0.18) | (±0.41) | (±0.25) | (±0.20) |

Table 3.7 Interface atom ratio (IR) and polarity ratio (PR) (with standard deviations in parentheses) for interfaces of PL complexes from the ABC, PIBASE, and Timbal datasets.

IR, interface atom ratio; PR, polarity ratio; PL, protein-ligand; Pi3, protein; Lj, ligand.

3.4 Conclusions

In this study, we characterized the residue and atom composition of overlapping protein-protein and protein-ligand interfaces from the ABC and PIBASE databases and compared these to a dataset derived from the Timbal database. According to the statistics, both interface types have, in general, a very similar composition. Among the three datasets of PL complexes, the protein interfaces of the Timbal dataset contain more hydrophobic residues and fewer polar residues than the two other datasets. The ligands in the PIBASE and Timbal datasets bind more flat on the protein surfaces or bind deeper into pockets on the protein surface than ABC ligands. Depending on the respective application in a ligand design project, researchers may consider to bias their principal dataset in one or the other direction. Selecting the appropriate set of reference data may slightly affect the physiochemical characteristics of designed ligands.

Chapter 4

STIM and ORAI Genes, Interactions with Transcription Factors, Differential Gene Expression and Co-expression Analysis on Breast Invasive Carcinoma Dataset

My contribution was to design the research project together with the co-authors Riccha Sethi, Mohamed Hamed, and Volkhard Helms. The analysis of transcription factors and STIM and ORAI genes was done by me and the differential gene expression, and co-expression analysis was done by Riccha Sethi. I and Volkhard Helms wrote the manuscript.

Abstract

Store-operated calcium (Ca²⁺) entry (SOCE) is ubiquitous mechanism for Ca²⁺ entry in eukaryotic cells, which regulates diverse cellular functions. SOCE is achieved primarily by the gating of the plasma membrane (PM)-localized-channel, ORAI, by the ER-localized Ca²⁺sensing protein, STIM. The discovery of transcription factor binding site (TFBS) motifs in specific locations on the STIM and ORAI promoters remains elusive. Moreover, the knowledge of the defects of STIM and ORAI genes expression and/or function linked to disease such as breast cancer is still obscured. Here, we used the HOCOMOCO and EPDnew databases to obtain a set of position weight matrix (PWM) and promoter sequences of STIM and ORAI, respectively, and mapped the possible binding motifs proteins using the FIMO tool. The results were then mapped with the set of transcription factors (TFs) targeting STIM and ORAI gene which were retrieved from the CheA database. We found ten predictive interactions between transcription factors bound to promoter regions of STIM and ORAI genes based on predictions using the STRING, prePPI, and mentha databases. Then, the collection of 63 TFs was used as gene of interest for co-expression and differential expression analysis on breast invasive carcinoma (BRCA) dataset. There, we found ORAI genes to be up-regulated, in contrast STIM1 and STIM2 were down-regulated. Unveiling the predicted transcription factors bound to the promoter regions of STIM and ORAI genes, the regulations of these genes and differential networks properties may suggest putative interactions for experimental studies and allows us to gain knowledge relation with breast cancer.

Keywords: STIM, ORAI, breast invasive carcinoma, promoters, transcription factors.

4.1 Introduction

Calcium (Ca²⁺) signals control many cellular functions ranging from short-term responses such as muscle contraction, impulse transmission and secretion to longer-term regulation of transcription, growth, and cell division. Store-operated Ca²⁺ entry (SOCE), the main Ca²⁺ influx mechanism in non-excitable cells, typically is activated in response to depletion of endoplasmic reticulum (ER)-Ca²⁺ stores [138,139]. SOCE is controlled by the ER-localized Ca²⁺ sensing proteins, STIM1 and STIM2 [140–142]. Recently, SOCE and the function of Ca²⁺ released-activated Ca²⁺ (CRAC) channels were shown to involve the ORAI1 or CRACM1 [143,144] and two homologs, ORAI2 and ORAI3 [145,146].

In mammals, the stromal interaction molecule (STIM) family has two members STIM1 (Figure 4.1) and STIM2 (Figure 4.2). They share a sequence identity almost 65%, they have diverse properties what results in different functions [147]. Both STIM are single spanning transmembrane (TM) proteins containing an N-terminal EF-hand domain responsible for calcium store sensing and a COOH terminal cytoplasmic domain [141,148–151]. STIM1 is a ubiquitously expressed, protein of 77 kDa and consists of 13 exons, located at chromosome 11, in 11p15.5. The STIM2 gene contains 148 amino acids residues (aa) longer than STIM1 (105-115 kDa) and also consists of 13 exons located at chromosome 4, in 4p15.1 [148].

The human ORAI family includes three members: ORAI1, ORAI2 and ORAI3 (Figure 4.3). They share a high sequence similarity of almost 89% and are broadly expressed with different expression levels depending on the cell type. These ORAI members are localized on the plasma membrane (PM) and consists of four transmembrane domains that are flanked by cytosolic NH₂ and COOH termini [152–155]. Many studies have shown that ORAI1 functions together with STIM1 to initiate CRAC currents [156,157]. On the other hand, previous studies showed that STIM1 and STIM2, are linked to several diseases such as Alzheimers disease (AD), Parkinson's and autoimmune diseases, brain and breast cancer, ischemia and, neurodegeneration diseases [158–164]. Moreover, a few studies found that ORAI1 and ORAI3 are linked to breast cancer migration and metastasis [165–168]. However, only few studies addressed the relation of STIM and ORAI genes to breast cancer.



Figure 4.1 Molecular structure of STIM1 (A) STIM1 domain organization. (B) Cartoon depicting a possible model of the STIM1 monomer in the resting state. The figure 4.1(B) was taken from [147].

N, N terminus; Ca²⁺, calcium; cEF, canonical EF-hand motif; nEF, non-canonical EF-hand motif; SAM, steril alpha motif; CC, coiled-coil domain; CAD, CRAC activation domain also called SOAR or CCb9; ER, endoplasmic reticulum; TM, transmembrane domain; S, serine-rich domain; S/P, serine- and proline-rich domain; K, poly-K and C, C terminus.



Figure 4.2 Molecular structure of STIM2 (A) STIM2 domain organization. (B) Cartoon depicting a possible model of the STIM2 monomer in the resting state. The figure 4.2(B) was taken from [147].

N, N terminus; Ca²⁺, calcium; cEF, canonical EF-hand motif; nEF, non-canonical EF-hand motif; SAM, steril alpha motif; CC, coiled-coil domain; CAD, CRAC activation domain also called SOAR or CCb9; ER, endoplasmic reticulum; TM, transmembrane domain; S, serine-rich domain; S/P, serine- and proline-rich domain; K, poly-K and C, C terminus.



Figure 4.3 Molecular structure of ORAI1, ORAI2, and ORAI3 (A) Domain organization of ORAI genes. (B) Cartoon depicting a possible model of ORAI in the resting state and sequence alignments. The figure 4.3(B) was taken from [169].

N, N terminus; P, Proline-rich; R, Arginine-rich; R/K, Arginine/Lysine-rich; TM, transmembrane domain; CC, coiled-coil domain and C, C terminus.

The binding of transcription factors (TFs) to specific locations in the genome is important for the coordination of transcriptional regulation in cells. Hence, the identification and accurate predictions of TF binding sites (TFBS) throughout genomes is a prerequisite to understanding of gene regulation and their networks [170]. Often, TFBSs are identified by using a set of target promoter sequences to characterize the binding properties for a particular TF that is known or assumed to bind these sequences. The position weight matrix (PWM) are computed from an alignment of these DNA sequences. Various databases are available which store TF targets such as the TRANSFAC [92] which contains experimentally verified TF target. There are also CheA [171] and Factorbook.org [172] which includes the resources of ChIP-seq and ChIP-chip data from the ENcyclopedia of DNA Element (ENCODE) project [173], HOCOMOCO [91], and JASPAR [90,174] databases. On the other hand, a number of online databases that store physical protein-protein interactions (PPIs) and networks have been developed such as the IntNetDB [175], mentha [176], OPHID [177], prePPI [178], Predictome [179], PIPs [180], and STRING [181,182]. Note that STRING also contains functional protein-protein (PP) interactions in additional to physical PPIs. These databases are the main resources of currently used for integrated research on PP interactions.

Here, by using publicly available databases and computational tools, we perform predict the transcription factors (TFs) responsible for the regulation of human STIM and ORAI genes, and perform differential gene expression and co-expression analysis for breast invasive carcinoma (BRCA) datasets. First, we created a dataset of transcription factors targeting STIM1, STIM2, ORAI1, ORAI2, and ORAI3 genes from the ChEA database and mapped them to the promoter regions of STIM and ORAI genes. Next, we searched for the predicted physical protein-protein interactions of TFs using the STRING, prePPI, and mentha databases. Then, the interacting TFs and bridge proteins were used for differential gene expression and co-expression analysis on datasets of level 3 RNASeqV2 breast invasive carcinoma (BRCA) from the TCGA database. Our main research question was to find the interacting transcription factors bound to the promoter regions of STIM1, STIM2, ORAI1, ORAI2, and ORAI3 and to assess how the regulation by these transcription factors may affect differential expression and co-expression analysis for breast invasive carcinoma datasets. The acquired knowledge should enhance our understanding of STIM and ORAI genes binding sites on promoters, their TF interactions and their regulations in differential expression gene and co-expression analysis.

4.2 Material and Methods

4.2.1 Transcription Factors Targeting STIM and ORAI Genes

We used the ChIP Enrichment Analysis (ChEA) database (amp.pharm.mssm.edu/chea) version 1 [171] to obtain the list of transcription factors (TF) targeted to STIM1, STIM2, and ORAI1-ORAI3 genes. We focused on transcription factors in human (*Homo sapiens*) only. The network of transcription factors targeting STIM1, STIM2, and ORAI1-ORAI3 genes was

visualized using the open-source platform for network analysis and visualization, Cytoscape (<u>http://www.cytoscape.org/</u>) [183].

4.2.2 Transcription Factors Binding Models

The dataset of transcription factor binding site (TFBS) models was obtained from the *Homo sapiens* Comprehensive Model Collection (HOCOMOCO) database version 9 (autosome.ru/HOCOMOCO/) [91]. We downloaded the probability matrices of AD curated collection in MEME text format (HOCOMOCOv9_AD_MEME.txt) which are derived against a uniform nucleotide background.

4.2.3 Sequence of Promoter Regions

A set of promoter regions of the STIM1, STIM2, ORAI1, ORAI2, and ORAI3 genes was downloaded on April 2016 from The Eukaryotic Promoter Database (EPDnew, epd.vitalit.ch/EPDnew_database.php) [89,184]. The EPDnew is a database of species-specific databases of experimentally validated promoters. The searched was filtered to *Homo sapiens*, and the default setting for the range from -499 to 100 bp relative to the transcription start site (TSS).

4.2.4 Motif Over-representation Analysis

A dataset of known binding motifs was obtained from the Find Individual Motif Occurrences (FIMO, meme-suite.org/tools/fimo) tool [185]. The FIMO tool is part of the Motif-Based Sequence Analysis Tools, the MEME Suite version 4.11.2 [186]. This tool searches a database of sequences for occurrences of known motifs and treats the motifs independently. The motif input from the HOCOMOCO database (HOCOMOCOv9_AD_MEME.txt) was uploaded in the "input the motifs" section of the FIMO tool. On the other hand, the promoter sequences of STIM1, STIM2, ORAI1, ORAI2, and ORAI3 were uploaded respectively each time in the "input the sequences" option.

4.2.5 Transcription Factors Predicted by CheA Mapped with the FIMO Results

We mapped 426 models from the HOCOMOCO database to our results of transcription factors targeted to STIM1, STIM2, and ORAI1-ORAI3 genes obtained from the ChEA database. Then the results were mapped according to their base pair (bp) location to each of the corresponding promoter regions.

4.2.6 Prediction of Physical Interactions of Protein-protein Complexes

Three different databases were used to obtain predicted physical interactions of protein-protein complexes, (i) the Search Tool for the Retrieval of Interacting Genes (STRING) database (<u>http://string-db.org</u>) downloaded in April 2016 [181,182], (ii) the prePPI database version 1.2.0 (<u>https://bhapp.c2b2.columbia.edu./PrePPI</u>) [178], and (iii) the mentha database (mentha.uniroma2.it/about.php) downloaded in April 2016 [176]. For the STRING database, the minimum required interaction score was set to the default cut-off of medium confidence (0.4). For the prePPI database, the probability \geq 0.5 indicates true interaction. Next, we grouped the interactions according to (i) direct interactions which contain no bridge protein and (ii) indirect interactions which require bridge protein/s.

4.2.7 TCGA BRCA Dataset

Publicly accessible RNA-Seq datasets version 2, level 3 of breast invasive carcinoma (BRCA) was downloaded in May 2016 from the online data portal, The Cancer Genome Atlas (TCGA; <u>http://tcga-data.nci.nih.gov</u>) [8]. We downloaded two separate datasets of 113 normal samples (n= 113) and 1102 tumor samples (n=1102), respectively. Each of them contains gene expression profiles. The datasets were pre-processed using python programming language and R software (<u>http://www.R-project.org</u>).

4.2.8 Computing Differential Gene Expression

We applied the widely used Bioconductor package DESeq2 version 1.12.3 [187] to identify genes showing differential expression in the TCGA dataset of breast invasive carcinoma. A false discovery rate (FDR) threshold of 0.1 after Benjamini and Hochberg multiple hypothesis correction was used for statistical significance genes. We used the pheatmap package to create a heatmap of the differentially expressed genes which used the normalized counts as the input.

4.2.9 Weighted Gene Co-expression Network Analysis

The co-expression analysis network was performed using the Weighted Gene Expression Network Analysis (WGCNA) [106] package version 1.51 for R version 3.2.0. We analyzed the results separately for normal and tumor samples. The clustering and principal-components analysis (PCA) was performed to identify outliers. In normal samples, 39 outlier samples were omitted. The cut-height was set to 8 and the soft-thresholding power parameter was set to 7. In tumor samples, 455 outlier samples were removed and the cut-height was set to 12 and the soft-thresholding power parameter was set to 4.

The significant modules obtained from normal and tumor samples were used to identify the differential interactions for all selected genes. The differential interactions were computed as following:

$$Differential Interactions = (normal - tumor) + (tumor - normal)$$
(4.1)

where (normal-tumor) are the interactions that are (edges) included in modules of the normal network which are not present in the tumor network modules and (tumor-normal) is vice versa. The edge weights greater than 0.02 obtained from the WGCNA were visualized by Cytoscape [183]. Edges with higher score are represented by thicker lines.

4.2.10 Flowchart

Figure 4.4 summarizes the workflow of the analysis of this study.



Figure 4.4 Flowchart summarizing the workflow of the analysis.

4.3 Results and Discussion

4.3.1 Transcription Factors Targeting STIM and ORAI Genes

Generally, high-throughput ChIP-seq data was used together with PWMs to identify the putative TF binding sites in the promoter regions of human STIM and ORAI genes. Thus, firstly we created a dataset of TFs targeting STIM and ORAI genes. These were downloaded from the ChEA database which contains gene lists of TFs obtained from ChIP-seq and ChIP-chip studies. We found a set of 23 TFs targeting STIM1, 29 (STIM2), 15 (ORAI1), 11 (ORAI2), and 13 (ORAI3), respectively. Overall, there are 48 non-redundant TFs targeting STIM and ORAI genes (Figure 4.5). Supplementary Information Table 4.1 lists the transcription factors targeted to STIM1, STIM2, and ORAI1-ORAI3, together with Pubmed ID (PMID), technique of experiments and cell types.



Figure 4.5 Dynamic graphical representations of 48 transcription factors targeting STIM and ORAI genes. The figure drawn using Cytoscape [183].

4.3.2 Promoter Sequences of STIM and ORAI Genes from the EPDnew Database

Two sets of promoter sequences each of STIM1 (named STIM1_1 and STIM1_2), and STIM2 (STIM2_1 and STIM2_2), and one promoter sequence for each ORAI1, ORAI2, and ORAI3 were retrieved from the EPDnew database. Supplementary Information Table 4.2 lists the promoter sequences for STIM1, STIM2, ORAI1, ORAI2, and ORAI3 genes. We used rather short promoter sequence region from -499 to 100 bp relative to the transcription start site (TSS) to focus on putative physical interactions between TFs targeted to these selected regions rather than longer promoter regions which may contain many unreliable interactions.

4.3.3 HOCOMOCO and FIMO Results

The STIM and ORAI promoter sequences were searched for occurrences of known motifs using the FIMO tool. This tool scans a sequence database for individual matches to each of the motifs provided in the HOCOMOCO dataset (HOCOMOCOv9_AD_MEME.txt) which contains 426 non-redundant curated binding models for 401 human TFs. In total, we found 411 known motifs in the promoter of STIM1_1, 79 for (STIM1_2), 276 for (STIM2_1), 83 for (STIM2_2), 304 for (ORAI1), 571 for (ORAI2), and 262 for (ORAI3), respectively.

Then, the respective motifs were mapped to the 91 transcription factors targeting STIM and ORAI genes obtained from the ChEA database. In total, there are 13 non-redundant transcription factors (E2F1, E2F4, EGR1, ELF1, ELK3, GATA1, HNF4A, MITF, MYC, PPARD, RUNX1, SOX2, and SPI1) that were found targeted STIM and ORAI genes. Table 4.1 lists the motifs found targeting STIM and ORAI genes. On the other hand, only one transcription factor was found in the promoters of STIM1_1, STIM1_2, and STIM2_2, respectively. Moreover, 12 TFs were found in the promoter of STIM2_1, 4 TFs in the promoters ORAI1 and ORAI2, and 5 TFs in the promoter ORAI3.

| Promoter Name | Motif | Strand | Start | End | n-value | a-value | Matched Sequence |
|------------------|----------|--------|-------|-----|-----------|----------|-------------------|
| STIM1 1 | | Stranu | 3tart | Enu | | | |
| STIM1 2 | ELK3_M | + | 494 | 505 | 4.050-05 | 0.0361 | GCCTGGAAGCCG |
| | SOX2_f1 | + | 443 | 458 | 2.14e-05 | 0.0238 | CTATGCATCAGAAAAG |
| 51111/2_1 | ELK3_f1 | + | 154 | 165 | 5.54e-05 | 0.0528 | CTCAGGATGTGG |
| | E2F1_f2 | - | 256 | 269 | 5.86e-06 | 0.00441 | AGGAGGCGGGGAAG |
| | E2F4_do | - | 256 | 269 | 4.41e-05 | 0.0372 | AGGAGGCGGGGAAG |
| | GATA1_si | - | 290 | 299 | 7.73e-05 | 0.0914 | GCTGATAACG |
| | EGR1_f2 | + | 339 | 349 | 9.25e-05 | 0.0115 | GGCGGGGCTGG |
| | EGR1_f2 | - | 354 | 364 | 8.56e-06 | 0.00306 | CGCGTGCGCGG |
| | E2F1_f2 | + | 404 | 417 | 8.14e-05 | 0.0204 | GGGAGGCGGGGGAT |
| | EGR1_f2 | + | 489 | 499 | 5.1e-06 | 0.00306 | GGAGGGGGCGG |
| | E2F1_f2 | + | 491 | 504 | 3.92e-05 | 0.0147 | AGGGGGCGGGGGGA |
| | EGR1_f2 | + | 510 | 520 | 1.94e-05 | 0.00361 | CGCGGCGGCGG |
| | EGR1_f2 | + | 513 | 523 | 1.23e-05 | 0.00306 | GGCGGCGGCGG |
| | EGR1_f2 | + | 516 | 526 | 7.86e-05 | 0.0115 | GGCGGCGGCGC |
| STIM2_2 | SOX2_f1 | - | 304 | 319 | 5.86e-05 | 0.0544 | TTTTACAAAATAATGA |
| ORAI1 | E2F4_do | + | 67 | 80 | 2.48E-005 | 0.0183 | AGTGGGCGCCAAAT |
| | E2F4_do | + | 280 | 291 | 6.29E-005 | 0.0231 | GGTGGGCGGGGAGC |
| | ELK3_f1 | - | 387 | 400 | 3.68E-005 | 0.0338 | TCCTGGAAGCGC |
| | PPARD_f1 | + | 396 | 409 | 3.91E-005 | 0.0446 | CGGGGCACAGGTGG |
| ORAI2 | RUNX1_f1 | + | 121 | 130 | 8.89E-005 | 0.101 | TCTGTGGGTA |
| | PPARD_f1 | + | 519 | 532 | 1.41E-005 | 0.0161 | TGGGCCACAGGCCA |
| | MYC_f1 | + | 528 | 538 | 2.47E-005 | 0.0114 | GGCCACGCGGC |
| | MYC_f1 | - | 529 | 539 | 5.66E-005 | 0.0131 | GGCCGCGTGGC |
| ORAI3 | HNF4A_f1 | + | 38 | 50 | 5.71e-06 | 0.0063 | GGACCAAAGGCCG |
| | MITF_f1 | + | 466 | 475 | 6.78e-05 | 0.0589 | ATCATGTGGC |
| | SPI1_si | - | 496 | 512 | 5.55e-07 | 0.000628 | CAAAACAGGAACTGGGA |
| | ELF1_f1 | - | 499 | 508 | 8.68e-05 | 0.0473 | ACAGGAACTG |
| | ELF1_f1 | - | 583 | 592 | 5.74e-05 | 0.0473 | CCAGGAAGAG |

Table 4.1 Motifs found targeting STIM and ORAI genes based on the CheA database and FIMO tool.

We prepared two diagrams to illustrate the transcription factors targeting STIM and ORAI on their promoter regions within the range of -499 to 100 bp, where 0 bp defines as the transcription start site (TSS) (Figure 4.6). We are assuming that possible physical interactions may occur between overlapping and neighbouring or adjacent transcription factors within \leq 50bp in the promoter regions. We identified nine such possible physical protein-protein interactions E2F1:E2F4, E2F1:EGR1, E2F1:GATA1, and E2F4:GATA1 in the promoter region of STIM2_1 and PPARD: E2F4 in the promoter region of ORAI1, PPARD: MYC (promoter region of ORAI2), MITF: SPI1, MITF: ELF1, and SPI1:ELF1 (promoter region of ORAI3) (Table 4.2). However, no putative interaction was found on promoter regions of

STIM1_1, STIM1_2 and STIM2_2 because only one transcription factor bound at each of them. In addition, an experimental study by Eylenstein et al. 2012 [188], found that NFKB1 are related to STIM1 and ORAI1. The FIMO tool predicted NFKB1 to targeted STIM1 and ORAI2, but these were not found in the CheA database. Though, we considered the interactions between NFKB1 and RUNX1 in the promoter region of ORA12.



Figure 4.6 Schematic illustration of transcription factors binding motifs in the promoters of the genes STIM1, STIM2, ORAI1, ORAI2, and ORAI3. (A) Transcription factor binding motifs in the promoters of genes STIM1 and STIM2. (B) Transcription factors binding motifs in the promoters of genes ORAI1, ORAI2, and ORAI3.

TSS, transcription start site; Promoter STIM1_1, first promoter of STIM1; Promoter STIM1_2, second promoter of STIM1; Promoter STIM2_1, first promoter of STIM2; Promoter STIM2_2, second promoter of STIM2; Promoter ORAI1, promoter of ORAI1; Promoter ORAI2, promoter of ORAI2; Promoter ORAI3, promoter of ORAI3.

Table 4.2 List of ten predicted physical transcription factor interactions, location in the promoter and their overlap or gap by base pair.

| | | | Location on p | promoter | |
|---------------|-------|-------|---------------|--------------|-------------------|
| Promoter Name | TF1 | TF2 | TF1 | TF2 | Overlap/gap by bp |
| STIM2_1 | E2F1 | E2F4 | -245 to -230 | | overlap by 15bp |
| | EGR1 | E2F1 | -145 to -135 | -95 to -82 | gap by 40bp |
| | | | -10 to 0 | -8 to 5 | overlap by 8bp |
| | *E2F1 | EGR1 | -8 to 5 | 11 to 21 | gap by 6bp |
| | | | -8 to 5 | 14 to 24 | gap by 9bp |
| | | | -8 to 5 | 17 to 27 | gap by 12bp |
| | E2F1 | GATA1 | -245 to -230 | -209 to -200 | gap by 21bp |
| | E2F4 | GATA1 | -245 to -230 | -209 to -200 | gap by 21bp |
| ORAI1 | PPARD | E2F4 | -112 to -99 | -103 to -90 | overlap by 4bp |
| ORAI2 | NFKB1 | RUNX1 | -430 to -420 | -378 to -369 | overlap by 42bp |
| | | | -397 to -387 | -378 to -369 | overlap by 9bp |
| | PPARD | MYC | 20 to 30 | 29 to 40 | overlap by 1bp |
| | | | 20 to 30 | 30 to 40 | gap by 0bp |
| ORAI3 | MITF | SPI1 | -33 to -24 | -3 to 13 | gap by 21bp |
| | MITF | ELF1 | -33 to -24 | 0 to 9 | gap by 24bp |
| | SPI1 | ELF1 | -3 to 13 | 0 to 9 | overlap by 9bp |

*The interaction between EGR1:E2F1 is assumed to be the same as interaction between E2F1:EGR1.

bp, base pair; TF1, first transcription factor and TF2, second transcription factor.

4.3.4 Physical Protein-protein Interactions

We used three well known databases STRING, prePPI, and mentha, to search for putative physical protein-protein interactions between transcription factors targeting STIM and ORAI genes on their promoter regions. Generally, these databases will search for possible physical protein-protein interactions either as direct or indirect interactions involving further bridge protein/s. Every database presented different results, scores, outputs, and networks.

Table 4.3 lists the results obtained from the STRING database. Figure 4.7 shows the network of interactions between 13 TFs targeting STIM and ORAI genes generated by the STRING database. The STRING database predicted eight possible protein-protein interactions either direct or involving bridge protein/s. The interactions of node1 and node2 which have a score of ≥ 0.5 suggest true interactions (Supplementary Information Table 4.3). In total, three direct interactions were observed between E2F1:E2F4, E2F4:GATA1 in the promoter of STIM2_1, and MITF:SPI1 in the promoter of ORAI3. Additionally, besides being direct interactions, E2F1:E2F4 was also connected by the bridge protein MYC and E2F4:GATA1 by MYC and EGR1.

Table 4.3 List of the predicted interactions between transcription factors obtained from the STRING database.

| Promoter name | Putative Interactions | Type of interaction | Bridge Proteins | Interactions |
|------------------|--------------------------|---------------------|--------------------|----------------------|
| STIM2_1 | E2F1:E2F4 | Direct | | |
| | E2F1:E2F4 | Bridge protein | MYC | E2F1>MYC>E2F4 |
| STIM2_1 | E2F1:EGR1 | Bridge protein | SPI1 | E2F1>SPI1>EGR1 |
| | E2F1:EGR1 | Bridge protein | E2F4 | E2F1>E2F4>EGR1 |
| | E2F1:EGR1 | Bridge protein | MYC | E2F1>MYC>EGR1 |
| STIM2_1 | E2F1:GATA1 | Bridge protein | ELF1 | E2F1>ELF1>GATA1 |
| | E2F1:GATA1 | Bridge protein | SPI1 | E2F1>SPI1>GATA1 |
| | E2F1:GATA1 | Bridge protein | MYC | E2F1>MYC>GATA1 |
| STIM2_1 | E2F4:GATA1 | Direct | | |
| | E2F4:GATA1 | Bridge protein | MYC | E2F4>MYC>GATA1 |
| | E2F4:GATA1 | Bridge protein | EGR1 | E2F4>EGR1>GATA1 |
| ORAI1 | PPARD:E2F4 | Bridge protein | EGR1 | PPARD>EGR1>E2F4 |
| ORAI2 | NFKB1:RUNX1 | Bridge protein | EGR1 | NFKB1>EGR1>RUNX1 |
| | NFKB1:RUNX1 | Bridge protein | MYC | NFKB1>MYC>RUNX1 |
| | NFKB1:RUNX1 | Bridge protein | SPI1 | NFKB1>SPI1>RUNX1 |
| ORAI2 | PPARD:MYC | Bridge protein | EGR1 | PPARD>EGR1>MYC |
| | PPARD:MYC | Bridge protein | HNF4A | PPARD>HNF4A>MYC |
| ORAI3 | MITF:ELF1 | Bridge proteins | SPI1, E2F1 | MITF>SPI1>E2F1>ELF1 |
| | MITF:ELF1 | Bridge proteins | SPI1, GATA1 | MITF>SPI1>GATA1>ELF1 |
| ORAI3 | MITF:SPI1 | Direct | | |
| ORAI3 | SPI1:ELF1 | Bridge protein | GATA1 | ELF1>GATA1>SPI1 |
| | SPI1:ELF1 | Bridge protein | E2F1 | ELF1>E2F1>SPI1 |



Figure 4.7 Network of 14 transcription factors targeting STIM and ORAI genes. Evidence view of the STRING database output depicting the transcription factors targeting STIM and ORAI genes obtained from <u>http://string-db.org/</u>.

The prePPI database results include prePPI LR, database LR, final probability and prediction code. The prediction code was labelled as S, T, G, E, M, C, and P to represent their sources of evidence used in the prediction. In our cases, the code of S was found which referred to structural modelling between the interactions of E2F1:E2F4, E2F1:GATA1, and NFKB1:RUNX1, T represents protein-peptide modelling (PPARD:E2F2 and PPARD:MYC) and G means GO term similarity (function similarity) (E2F1:EGR1 and MITF:SPI1). Generally, the prePPI database predicts only direct interactions. In total, four putative physical protein-protein interactions are predicted true by the prePPI database with a final probability \geq 0.5. This includes the interactions between E2F1:E2F4 (final probability of 1.00), E2F:EGR1 (final probability of 0.68), MITF:SPI1 (final probability of 1.00), and SPI1:ELF1 (final probability of 0.99), respectively (Table 4.4).

Table 4.4 List of predicted interactions between transcription factors obtained from the prePPI database.

| Promoter name | Predicted Interactions | prePPI LR | Database LR | Final Probability | Prediction Code |
|---------------|---------------------------|--------------|---------------|----------------------|--------------------|
| STIM2_1 | E2F1:E2F4 | 109392.00 | 957.82 | 1.00 | S |
| STIM2_1 | E2F1:EGR1 | 1268.49 | Not available | 0.68 | G |
| STIM2_1 | E2F1:GATA1 | 31.17 | Not available | 0.05 | S |
| STIM2_1 | E2F4:GATA1 | Not found | Not found | Not found | Not found |
| ORAI1 | PPARD:E2F4 | 23.76 | Not available | 0.04 | Т |
| ORAI2 | NFKB1:RUNX1 | 388.29 | Not available | 0.39 | S |
| ORAI2 | PPARD:MYC | 61.64 | Not available | 0.09 | Т |
| ORAI3 | MITF:ELF1 | Not found | Not found | Not found | Not found |
| ORAI3 | MITF:SPI1 | 47.73 | 4625.64 | 1.00 | G |
| ORAI3 | SPI1:ELF1 | 60755.6 | Not available | 0.99 | S |

LR, likelihood ratio; G, GO term similarity; S, structural modelling and T, protein-peptide modelling.

Generally, the mentha database can predict direct and indirect interactions. Table 4.5 lists the predicted interactions obtained from the mentha database. In total, three interactions with bridge proteins are predicted between E2F1:EGR1, PPARD:E2F4, PPARD:MYC, and one direct interaction between E2F1:E2F4. We identified one bridge protein (HDAC1) in the NFKB1:RUNX1 interaction, three bridge proteins (HDAC1, HDAC3, and NCOR2) in PPARD:E2F4 interactions, four bridge proteins (CDKN2A, CREBBP, EP300, and SP1) associate with E2F1:EGR1 interactions, and five bridge proteins (EP300, HDAC1, HDAC2, HDAC3, and KDM1A) with PPARD:MYC interactions, respectively. However, interactions

between E2F1:GATA1, E2F4:GATA1, MITF:SPI1, MITF:ELF1, and SPI1:ELF1 were not found in mentha.

| Promoter name | Putative Interactions | Type of interaction | Bridge Proteins | TF1 | TF2 | Score | PMID |
|------------------|--------------------------|---------------------|--------------------|--------|-------|-------|---|
| STIM2_1 | E2F1:E2F4 | Direct | NULL | E2F4 | E2F1 | 0.126 | 16357170 |
| STIM2_1 | E2F1:EGR1 | Bridge protein | CDKN2A | CDKN2A | E2F1 | 0.507 | 11314038 |
| | | | | CDKN2A | EGR1 | 0.623 | 19057511 |
| | E2F1:EGR1 | Bridge protein | SP1 | E2F1 | SP1 | 0.902 | 10547281, 8657141, 10409740, 8657142 |
| | | | | EGR1 | SP1 | 0.523 | 20121949 |
| | E2F1:EGR1 | Bridge protein | CREBBP | CREBBP | E2F1 | 0.523 | 12748276, 8932363 |
| | | | | CREBBP | EGR1 | 0.623 | 9806899 |
| | E2F1:EGR1 | Bridge protein | EP300 | EP300 | E2F1 | 0.507 | 24112038, 15123636, 23001041 |
| | | | | EP300 | EGR1 | 0.93 | 9806899, 15225550, 20089040, 20018936 |
| STIM2_1 | E2F1:GATA1 | | | | | | |
| | Not found | | | | | | |
| STIM2_1 | E2F4:GATA1 | | | | | | |
| | Not found | | | | | | |
| ORAI1 | PPARD:E2F4 | Bridge protein | NCOR2 | NCOR2 | PPARD | 0.454 | 11867749 |
| | | | | NCOR2 | E2F4 | 0.376 | 22508987 |
| | PPARD:E2F4 | Bridge protein | HDAC1 | HDAC1 | PPARD | 0.569 | 18037904, 11867749 |
| | | | | E2F4 | HDAC1 | 0.523 | 9724731, 23060449, 9858615 |
| | PPARD:E2F4 | Bridge protein | HDAC3 | PPARD | HDAC3 | 0.454 | 11867749,12943985 |
| | | | | E2F4 | HDAC3 | 0.376 | 22508987 |
| ORAI2 | NFKB1:RUNX1 | Bridge protein | HDAC1 | NFKB1 | HDAC1 | 0.91 | 25241761 16319923 24448807 12972430 17827154 11931769 17962807 |
| | | | | HDAC1 | RUNX1 | 0.73 | 22498736, 16652147, 21059642 |
| ORAI2 | PPARD:MYC | Bridge protein | HDAC3 | PPARD | HDAC3 | 0.454 | 11867749, 12943985 |
| | | | | MYC | HDAC3 | 0.692 | 22002311,18483244, 23079660 |
| | PPARD:MYC | Bridge protein | KDM1A | PPARD | KDM1A | 0.49 | 23455924 |
| | | | | MYC | KDM1A | 0.332 | 23455924 |
| | PPARD:MYC | Bridge protein | HDAC2 | HDAC2 | PPARD | 0.309 | 25241761, 11867749 |
| | | | | HDAC2 | MYC | 0.472 | 17314511, 20195357, 22286234 |
| | PPARD:MYC | Bridge protein | HDAC1 | HDAC1 | PPARD | 0.569 | 18037904, 11867749 |
| | | | | HDAC1 | MYC | 0.911 | 22286234, 18003922, 26496610, 17314511, 18271930, 24951594 |

Table 4.5 List of the predicted interactions between transcription factors obtained from the mentha database.

| | PPARD:MYC | Bridge protein | EP300 | EP300 | PPARD | 0.309 | 16930961 |
|-------|-----------|----------------|-------|-------|-------|-------|---|
| | | | | EP300 | MYC | 0.911 | 17157259, 15616592, 16287840, 16126174 |
| ORAI3 | MITF:SPI1 | | | | | | |
| | Not found | | | | | | |
| ORAI3 | MITF:ELF1 | | | | | | |
| | Not found | | | | | | |
| ORAI3 | SPI1:ELF1 | | | | | | |
| | Not found | | | | | | |

PMID, Pub Med ID.

Table 4.6 summarizes the type of interactions either as direct interaction or using bridge protein/s predicted from the STRING, prePPI, and mentha databases. If one of these databases predicts any interaction; which the first TF (TF1) interact with the second TF (TF2), we labelled "YES" in "Consider as TFs pair?" column. Overall, 14 non-redundant TFs are found to act as bridge proteins predicted by the STRING and mentha databases. Next, we used these 63 non-redundant genes (including 48 TFs predicted by the ChEA database and STIM1, STIM2, ORAI1, ORAI2, and ORAI3 genes itself) as genes of interest for differential expression gene analysis and co-expression analysis of the breast invasive carcinoma dataset (Table 4.7).

| Promoter | <u>1 1 1, una </u> | inonina a | Type of | | | | Consider |
|----------|--------------------|-----------|----------------|--------|--------|--------|-------------|
| name | TF1 | TF2 | Interaction | STRING | prePPI | mentha | as TF pair? |
| STIM2_1 | E2F1 | E2F4 | Direct | YES | YES | YES | YES |
| | | | Bridge Protein | YES | | | |
| | E2F1 | EGR1 | Direct | | YES | | YES |
| | | | Bridge Protein | YES | | YES | |
| | E2F1 | GATA1 | Direct | | | | |
| | | | Bridge Protein | YES | | | YES |
| | E2F4 | GATA1 | Direct | YES | | | YES |
| | | | Bridge Protein | YES | | | |
| ORAI1 | PPARD | E2F4 | Direct | | | | |
| | | | Bridge Protein | YES | | YES | YES |
| ORAI2 | NFKB1 | RUNX1 | Direct | | | | |
| | | | Bridge Protein | YES | | YES | YES |
| | PPARD | MYC | Direct | | | | |
| | | | Bridge Protein | YES | | YES | YES |
| ORAI3 | MITF | ELF1 | Direct | | | | |
| | | | Bridge Protein | YES | | | YES |
| | MITF | SPI1 | Direct | YES | YES | | YES |
| | | | Bridge Protein | | | | |
| | SPI1 | ELF1 | Direct | | YES | | YES |
| | | | Bridge Protein | YES | | | |

Table 4.6 List of predicted interactions between transcription factors obtained from the STRING, prePPI, and mentha databases.

| AR | FOXA1 | HNF4A | PAX3-FKHR | STIM1 |
|--------|-------|--------|-----------|--------|
| BCL6 | FOXA2 | HOXB7 | PHF8 | STIM2 |
| CDKN2A | FOXP1 | KDM1A | POU3F2 | TFAP2C |
| CREBBP | FOXP2 | KLF5 | PPARD | TFEB |
| CUX1 | GABP | MITF | RBPJ | TOP2B |
| E2F1 | GATA1 | MYC | RUNX1 | TP63 |
| E2F4 | GATA2 | MYCN | RUNX2 | TRIM28 |
| EBNA2 | GATA3 | NCOR1 | SCL | TTF2 |
| EGR1 | GATA4 | NCOR2 | SMAD4 | VDR |
| ELF1 | GATA6 | *NFKB1 | SOX11 | WT1 |
| ELK3 | HDAC1 | ORAI1 | SOX2 | ZNF217 |
| EP300 | HDAC2 | ORAI2 | SP1 | |
| FLI1 | HDAC3 | ORAI3 | SPI1 | |

Table 4.7 List of 63 non-redundant genes of interest obtained from the CheA, STRING, prePPI, and mentha databases.

* NFKB1 was found to be related to STIM and ORAI genes by experimental study and was predicted by the FIMO tool.

4.3.5 Differential Gene Expression Analysis

To identify specific genes that were differentially expressed in breast invasive carcinoma dataset, we used DESeq2 packages. This dataset contains 113 normal and 1102 tumor samples. Aforementioned, we used a set of 63 genes of interest (Table 4.7) for differentially expressed analysis with STIM and ORAI genes as the main focus. By setting the p-value ≤ 0.05 and FDR value of 0.1, DESeq2 analysis identified 45 out of 63 genes that were differentially expressed with 26 genes being significantly up-regulated and 19 genes are down-regulated (Table 4.8). The ten most significantly differentially up-regulated genes were E2F1, ORAI2, CDKN2A, SOXII, GATA3, TRIM28, RUNX2, FOXA1, HDAC1, and KDM1A. On the other hand, the ten most significantly down-regulated genes were EGR1, FLI1, FOXP2, BCL6, TP63, MITF, SMAD4, STIM2, RBPJ, and ELK3. The ORAI1, ORAI2, and ORAI3 genes were found up-regulated. This agrees with the results of several studies stating that the expression of ORAI1 [167,189,190] and ORAI3 [168,190] genes increased in primary human breast cancer cells or tissues. On the other hand, we found STIM1 and STIM2 to be down-regulated genes.

Following data pre-processing, a total of 59 genes were identified. Four genes EBNA2, GABP, PAX3-FKHR, and SCL were not found in both normal and tumor samples of the breast invasive carcinoma dataset. Based on principal component analysis (PCA) for the gene expression profile data, it was found that the gene expression level in tumor (breast invasive carcinoma) samples was dissimilar to that in normal samples, thus the two groups were distinguished absolutely at the gene expression level (Figure 4.8).

| | Log fold | | FDR- | | Log fold | | FDR- | |
|------------------------|----------|-----------|-----------|--------------------------|----------|-----------|-----------|--|
| | change | p-value | adjusted | | change | p-value | adjusted | |
| (a) Up-regulated genes | | | | (b) Down-regulated genes | | | | |
| E2F1 | 2.81 | 5.29E-132 | 2.81E-130 | EGR1 | -2.68 | 7.75E-84 | 2.05E-82 | |
| ORAI2 | 1.05 | 2.78E-63 | 4.92E-62 | FLI1 | -1.22 | 3.04E-55 | 3.22E-54 | |
| CDKN2A | 2.95 | 8.34E-63 | 1.10E-61 | FOXP2 | -2.38 | 5.07E-50 | 4.48E-49 | |
| SOX11 | 3.06 | 2.24E-43 | 1.19E-42 | BCL6 | -1.12 | 2.54E-49 | 1.92E-48 | |
| GATA3 | 1.87 | 5.02E-42 | 2.42E-41 | TP63 | -2.56 | 6.32E-45 | 4.19E-44 | |
| TRIM28 | 0.83 | 3.19E-40 | 1.41E-39 | MITF | -1.21 | 1.55E-44 | 9.12E-44 | |
| RUNX2 | 1.28 | 6.42E-38 | 2.62E-37 | SMAD4 | -0.53 | 1.13E-28 | 4.00E-28 | |
| FOXA1 | 1.86 | 3.20E-33 | 1.21E-32 | STIM2 | -0.50 | 1.56E-28 | 5.16E-28 | |
| HDAC1 | 0.59 | 1.51E-27 | 4.70E-27 | RBPJ | -0.49 | 2.32E-25 | 6.83E-25 | |
| KDM1A | 0.52 | 1.80E-22 | 5.02E-22 | ELK3 | -0.69 | 2.40E-21 | 6.07E-21 | |
| MYCN | 1.75 | 8.31E-22 | 2.20E-21 | MYC | -0.90 | 7.03E-16 | 1.49E-15 | |
| ZNF217 | 0.72 | 1.02E-18 | 2.45E-18 | SP1 | -0.25 | 4.66E-12 | 9.14E-12 | |
| HDAC2 | 0.61 | 6.77E-17 | 1.56E-16 | KLF5 | -0.82 | 4.02E-08 | 6.45E-08 | |
| TTF2 | 0.57 | 2.97E-16 | 6.55E-16 | NCOR1 | -0.31 | 1.637E-06 | 2.55E-06 | |
| ORAI1 | 0.55 | 4.46E-13 | 9.10E-13 | TFAP2C | -0.47 | 2.00E-05 | 2.87E-05 | |
| SPI1 | 0.67 | 7.14E-12 | 1.35E-11 | STIM1 | -0.26 | 3.01E-05 | 4.10E-05 | |
| VDR | 0.47 | 2.44E-11 | 4.45E-11 | EP300 | -0.18 | 0.0009369 | 0.0011822 | |
| PHF8 | 0.36 | 3.40E-11 | 6.01E-11 | GATA6 | -0.33 | 0.006 | 0.007 | |
| PPARD | 0.39 | 1.52E-10 | 2.59E-10 | FOXP1 | -0.17 | 0.013 | 0.015 | |
| ORAI3 | 0.46 | 1E-09 | 1.69E-09 | | | | | |
| CUX1 | 0.30 | 3.87E-06 | 5.86E-06 | | | | | |
| GATA2 | 0.68 | 5.90E-06 | 8.69E-06 | | | | | |
| NCOR2 | 0.22 | 2.66E-05 | 3.71E-05 | | | | | |
| HDAC3 | 0.17 | 3.45E-05 | 4.57E-05 | | | | | |
| RUNX1 | 0.25 | 0.0004872 | 0.0006298 | | | | | |
| TOP2B | 0.11 | 0.023 | 0.027 | | | | | |

Table 4.8 List of differentially expressed genes in the dataset of breast invasive carcinoma (BRCA).

FDR, False Discovery Rate.


Figure 4.8 Plot of principal component analysis of 113 normal samples and 1102 tumor samples.

PC, principal component.

4.3.6 Results of Gene Co-expression Analysis

The analysis of gene co-expression of 113 normal samples and 1102 tumor samples was done separately using the WGCNA program.

4.3.6.1 Normal Samples

By observing the result obtained from the clustering dendogram and the topological overlap matrix (TOM) heatmap for all genes in normal samples (Figure 4.9), we found two significant modules. The blue module contains 15 genes (including ORAI2) and the turquoise module contains 18 genes (Table 4.9).



Figure 4.9 Topological overlap matrix (TOM) heatmap corresponding to the two coexpression modules in normal samples. Each row and column of the heatmap represents a single gene. Red indicates high levels of co-expression genes. The dendograms on the upper and left sides show the hierarchical clustering tree of genes.

4.3.6.2 Tumor Samples

In the tumor samples, we identified three significant modules. The brown module contains eight genes, the blue module 11 genes, and the turquoise module 13 genes, respectively (Figure 4.10). On the other hand, we found STM2 in brown module and ORAI1 in blue module. Table 4.9 summarizes all genes found in normal and tumor samples. The expression patterns of the individual modules will provide some clue about their functions.



Figure 4.10 Topological overlap matrix (TOM) heatmap corresponding to the three coexpression modules in tumor samples. Each row and column of the heatmap represents a single gene. Red indicates high levels of co-expression genes. The dendograms on the upper and left sides show the hierarchical clustering tree of genes.

Table 4.9 Number of genes and gene name of the significant modules obtained from the normal and tumor samples identified by the WGCNA program.

| Module | No. of genes | Gene Na | me | | Module | No. of genes | Gene Na | ime | |
|--------------------|-----------------|---------|--------|-------------------|-----------|-----------------|---------|--------|--------|
| (a) Normal samples | | | | (b) Tumor samples | | | | | |
| Blue | 15 | FOXP1 | TOP2B | CUX1 | Blue | 11 | SMAD4 | PHF8 | CREBBP |
| | | SMAD4 | ORAI2 | AR | | | ELF1 | SP1 | |
| | | FOXA1 | PHF8 | EP300 | | | NCOR1 | CDKN2A | |
| | | TTF2 | SP1 | GATA3 | | | TOP2B | AR | |
| | | NCOR1 | CDKN2A | CREBBP | | | ORAI1 | EP300 | |
| Turquoise | 18 | FOXP2 | NCOR2 | HOXB7 | Turquoise | 13 | FOXP2 | NFKB1 | RUNX1 |
| | | HDAC1 | TP63 | VDR | | | HDAC1 | TRIM28 | |
| | | HDAC3 | RBPJ | GATA6 | | | MITF | GATA6 | |
| | | KDM1A | KLF5 | GATA2 | | | TP63 | E2F1 | |
| | | MITF | TRIM28 | ELK3 | | | EGR1 | RUNX2 | |
| | | HDAC2 | FLI1 | TFAP2C | | | KLF5 | ELK3 | |
| | | | | | Brown | 8 | FOXA1 | ZNF217 | GATA3 |
| | | | | | | | STIM2 | FLI1 | E2F4 |
| | | | | | | | SPI1 | TFEB | |

4.3.7 Differential Interaction Networks

Following the analysis with the WGCNA package, we extracted all the genes differentially expressed from the significant modules of the TOM plot to construct networks of differential interactions using the Cytoscape [183]. The differential interactions were computed to observe which interactions are not in normal samples (labelled as normal-tumor) and vice versa for tumor samples (tumor-normal). We used the edges scores of the threshold ≥ 0.02 . This resulted in 83 differential interactions in normal-tumor samples and 61 interactions in tumor-normal samples (Supplementary Information Table 4.4). The differential interactions networks suggesting that one of our focus genes, ORAI2 interacts with the SP1 and SMAD4 genes in normal-tumor, while STIM2 interacts with the ELK3, FOXP2 and FLI1 in tumor-normal (Figure 4.11). In the normal-tumor network, the three highest scored edges are the interactions between EP300:CREBPP (0.2822), NCOR1:EP300 (0.2343), and SP1:EP300 (0.2093), respectively. In the tumor-normal network, the three highest scored edges are the interactions between FOXP2:ELK3 (0.1546), MITF:RUNX2 (0.1528), and FOXA1:GATA3 (0.1112), respectively. On the other hand, genes with the highest number of edges which known as hubs play centred roles in the analysis of the networks. Thus, we identified two genes CREBBP and ELK3 as hub genes in the normal-tumor network which contains eight edges (Figure 4.11 (A)) while in the tumor-normal network, ELK3 is the hub gene which contains 12 edges (Figure 4.11 (B)).



Figure 4.11 Differential Interaction Networks (A) Differential interactions network of normaltumor. (B) Differential interactions network of tumor-normal.

Red edges indicate interactions that found in differential interactions and black edges indicate the interactions which were not found in differential interactions. Blue nodes indicate the genes found in the blue module, brown (blue module), and turquoise (turquoise module), respectively.

4.3.8 Regulation of STIM and ORAI Genes on Normal and Tumor Breast Cancer Tissue As mentioned above, we identified ten putative interactions of transcription factors bound to promoter region of STIM and ORAI. There, five of the interactions involving E2F1 and E2F4 on promoter regions of STIM2 and ORAI1 (Table 4.6). Previous studies showed that the transcription factors of E2F family (E2F1-8) is recognized to regulate many important genes and involved in many biological processes such as apoptosis, cell proliferation, differentiation, and DNA damage response [191,192]. Furthermore, we identified E2F1 was up-regulated gene in tumor samples and interact with RUNX1 in differential interaction of tumor-normal network. Moreover, our results show that E2F1 and E2F4 are found co-expressed in tumor samples. We speculate that E2F1 and E2F4 are regulating STIM2 and ORAI1 genes which also show relation to breast cancer. In addition, we identified an interaction of NFKB1:RUNX1 on the promoter region of ORAI2 which also found co-expressed in turquoise module of tumor samples. NFKB1 is one of the family members of NFKB1 transcription factors [193]. Several studies noted that NFKB among the Ca²⁺ sensitive transcription factors which are associate with STIM1 and ORAI1 to stimulate cell proliferation and differentiation [194-196]. On the other hand, two interactions are found involving the micophthalmia-associated transcription factor (MITF) on promoter region of ORAI3 (MITF:ELF1 and MITF:SPI1). We found that MITF are co-expression in both turquoise modules of normal and tumor samples. A study by Carmit and co-workers noted that MITF functions as master regulatory of melanocytes development and melanoma oncogene [197]. Furthermore, Stanisz and co-workers stated the role of STIM and ORAI in melanocytes and melanoma which significantly correlates with the expression of MITF [163,198]. Generally, our findings suggest several roles of STIM and ORAI genes in normal and breast cancer tissues.

4.4 Conclusions

In this work, we identified 13 non-redundant transcription factors targeting STIM and ORAI genes. We then found ten putative TFs interactions bound on promoter regions of STIM and ORAI genes predicted by the STRING, prePPI, and mentha databases. According to these interactions, we found 14 non-redundant TFs act as bridge proteins. Then, we identified 63 non-redundant genes as genes of interest for differential expression gene analysis and co-expression analysis of breast invasive carcinoma dataset. In the differential expression analysis, we found 26 up-regulated genes including ORAI1, ORAI2, and ORAI3, while the 19 down-regulated genes include STIM1 and STIM2. On the other hand, in the co-expression analysis,

we found two significant modules (blue and turquoise modules) in normal samples and three significant modules (brown, blue, and turquoise modules) in tumor samples which the expression patterns of the individual modules tend to provide clues about their functions. Next, we identified 83 differential interactions in normal-tumor samples and 61 interactions in tumor-normal samples. Finally, we identified CREBBP and ELK3 as hubs genes in the normal-tumor network, and ELK3 in the tumor-normal network. Overall, our findings form an important basis for identifying TFs targeting STIM and ORAI genes and demonstrate the significant involvements of STIM and ORAI genes in breast cancer. In ongoing work, we are extending the gene enrichment analysis and applying the framework to datasets of diseases reported linked to STIM and ORAI genes such as Alzheimers disease (AD) and prostate cancer.

Chapter 5

Evaluation of the Protein Pocket Identification Tools on Protein-Ligand Complexes

My contribution was to write the manuscript, designed the research project, and analyzed the results together with the co-authors Zhao Yuan, Rahmad Akbar, and Volkhard Helms. I and Rahmad Akbar co-supervised Zhao Yuan. Zhao Yuan performed the calculations.

Abstract

Binding pockets are regions on protein surfaces where substrates of enzymatic reactions or effector molecules and co-factors may bind. Thus, identifying these cavities is often a prerequisite step for structure-based drug design. Various computational methods have been developed to identify such sites on protein surfaces. In this work, we evaluated the seven tools DEPTH, DoGSiteScorer, Fpocket, GHECOM, IsoMif, PocketPicker, and ProACT2 on a dataset of 167 non-redundant protein-ligand complexes. We analyzed how well the predicted pocket-lining residues overlap with the residues that contact the ligand. We used the residue overlap to define a score as a measure of the predictive capabilities of the tools. Even though the tools predicted pockets of various sizes and shapes we found comparable performance amongst the predictions of five tools (DEPTH, GHECOM, DoGSiteScorer, Fpocket, and IsoMif) in terms of average score. Using always the most suitable tool improved the average score by 28% over randomly selecting a tool. To support users in a pocket prediction scenario, we trained a random forest model (classifier) to output a list of suitable tools for a given protein structure. This classifier should be useful for prioritizing the tools to be used for unknown proteins or proteins that are not contained in our dataset.

Keywords: classifier, protein-ligand complexes, protein pockets, and random forest.

5.1 Introduction

Proteins play major roles in practically all cellular processes. They typically interact with small molecules (ligands), nucleic acids or other proteins to perform a certain function. These interactions often occur in a particular site on the protein surface (binding site). As binding sites are thus often related directly to protein function, it is important to advance our understanding on these sites. The large collection of experimentally determined three-dimensional structures of protein-ligand complexes stored in the protein data bank (PDB) allows us to study these binding sites. For instance, it has been shown that binding sites for small molecule ligands tend to be rather hydrophobic with few selected polar and charged residues [199–201] and tend to be found in deep pockets on the proteins surface [53,202].

Based on this data, one can develop algorithms to identify cavities that may accommodate bound ligands on protein surfaces. The current batch of such algorithms fall into five categories (i) geometric methods that can be further grouped into the three subcategories grid system scanning, probe sphere filling, and alpha shape [49,50], (ii) energy based methods, (iii) evolution based methods, (iv) blind docking and molecular dynamics and, (v) combined approaches [48]. Grid system scanning basically projects a protein structure onto a threedimensional grid of points and examines spatial overlaps on this grid. DEPTH [57,58], DoGSiteScorer [56], GHECOM [49], and PocketPicker [55] are tools that implement a grid system scanning approach. Probe sphere filling methods generate a set of probe spheres to fill cavities on protein surfaces. Pockets are then defined as those regions containing the highest amount of spheres, e.g. by the tool IsoMif [63]. Alpha shape methods rely on the alpha-shape theory and Voronoi tessellation to identify a pocket. Fpocket [68] is a representative of alpha shape methods. Energy based methods identify pockets using energetic criteria. Cavities with the largest total interaction energies are defined as pockets. For instance, ProACT2 [69,70] is a representative of energy based method. Other tools employ further strategies. For example, Rate4Site [203] uses an evolution based approach and MolSite [77] utilizes blind docking. As a wide range of different strategies and approaches are employed by current pocket identification tools, it is of interest to compare and contrast the performance of these tools.

Defining the correct pockets on protein surfaces is not an easy task [68,204]. Schematically, Figure 5.1 sketches three possible ways to define the pocket volume of a pacman-shape surface cavity that is shown in two dimensions. It is unclear what definition is correct and most useful. Here we analyzed how well the constructed pockets overlap with the protein contacts made by small molecule ligands in their X-ray conformations. One should add,

as a word of caution, that native or synthetic ligands may either be smaller than surface pockets or exceed the volume of the pocket into the solution. Figure 5.2 shows for a case system that these tools predict pockets of various sizes and shapes. Indeed, in some cases the ligand is not fully enclosed by the detected pocket whereas other tools generated rather large pockets.



Figure 5.1 Two dimensional illustrations of three possible ways to define the pocket volume of a pacman-shape surface cavity.



Figure 5.2 Pockets (blue) identified by DEPTH (lining residues) (top panel left), GHECOM (top panel middle), Fpocket (top panel right), DoGSiteScorer (middle panel left), PocketPicker (middle panel middle), IsoMif (middle panel right), and ProACT2 (bottom panel) that overlap with the ligand HNT (red sticks) bound to human phenylethanolamine N-methyltransferase, PNMT (PDB ID: 2G70). The figures were generated using PyMOL Molecular Graphics System [47].

Here, we compared the performance of seven tools using a set of quality metrics on a set of 167 protein-ligand complexes. We then computed a set of physico-chemical and geometric features for each ligand-bound pocket in the dataset. Correlation analysis of pocket features and the quality metrics revealed only weak correlation between pocket features and the predictive performance of the tools. In general, we found comparable performance in the predictions of five tools DEPTH [57,58], GHECOM [49], DoGSiteScorer [56], Fpocket [68], and IsoMif [63].

5.2 Material and Methods

5.2.1 Dataset

We used a non-redundant dataset of 195 protein-ligand complexes retrieved by Degac et. al. [205] from the PDBbind [37,206] database version 2014. Beside the protein, these complexes contain a single ligand with weight less than 1000 Da and the resolution of the X-ray structure must be equal to or better than 2.5 Å. Additionally, the ligand molecules must contain only common organic elements and the protein molecules include only the standard 20 amino acid residues in the area of the binding sites. Complexes with insertions and/or residues numbered with special characters (20 complexes), or that could not be processed with Fpocket [68] (five complexes) or with DoGSiteScorer [56] (three complexes) were removed. This yielded a final dataset of 167 complexes (Supplementary Information Table 5.1). Only a single chain was considered if a complex contains more than one homomer.

5.2.2 Tools

Initially, we considered a total of 25 pocket identification tools. However, among these, we were only able to automate the use of seven tools due to various limitations. Table 1 lists the tools used in this work. The complete set of tools and the respective limitations are listed in Supplementary Information Table 5.2.

| Program | URL | Year |
|---------------|--|------|
| DEPTH | http://mspc.bii.a-star.edu.sg/tankp/help.html | 2011 |
| DoGSiteScorer | http://dogsite.zbh.uni-hamburg.de/ | 2012 |
| Fpocket | http://Fpocket.sourceforge.net/ | 2009 |
| GHECOM | http://strcomp.protein.osaka-u.ac.jp/GHECOM/ | 2010 |
| IsoMif | http://bcb.med.usherbrooke.ca/imfi.php | 2015 |
| PocketPicker | http://gecco.org.chemie.uni-frankfurt.de/pocketpicker/index.html | 2007 |
| ProACT2 | http://people.cryst.bbk.ac.uk/~ubcg66a/proact2_summary.html | 2010 |

Table 5.1 Names, URLs and year of creation of the seven protein pocket prediction tools.

5.2.3 Binding Site Definition

A binding site is defined as the set of residues which are located within 5.0 Å from the ligand surface (actual positive class in the X-ray structure of the protein-ligand complexes). The remaining residues in the protein are defined as non-binding site residues (actual negative class).

5.2.4 Model Evaluations

Predicted pockets were evaluated by comparing their pocket lining residues with the residues of the corresponding binding site of the respective protein-ligand complexes using a confusion matrix. Each column of the confusion matrix represents the predicted class, whereas each row represents the actual class. We used this matrix to quantify the correct and incorrect predictions of each tool. In our case, TP (true positive) is the overlap between residues in the binding site and the lining residues of the predicted pockets, FN (false negative) are the residues in the binding site residues that were predicted as non-binding site. FP (false positive) are the non-binding site residues that were correctly predicted as non-binding site. From the confusion matrix we computed the following quality metrics.

5.2.5 MCC

The Matthews correlation coefficient (MCC) [207] quantifies the degree of correlation between the actual and predicted classes of the residues. An MCC value of 1 indicates that all predictions are correct, -1 for completely incorrect predictions [208], and a value of zero indicates a random prediction [209]. MCC is defined in equation 5.1:

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FN)(TP + FP)(TN + FP)(TN + FN)}}$$
(5.1)

5.2.6 Precision

Precision (P) is the proportion of correct predictions (TP) among all positive predictions (TP and FP). The precision value ranges from 0 to 1.

$$P = \frac{TP}{TP + FP} \tag{5.2}$$

5.2.7 Recall

Recall (R) is the proportion of correct predictions (TP) in condition positive (TP and FN).

The recall value ranges from 0 to 1.

$$R = \frac{TP}{TP + FN} \tag{5.3}$$

5.2.8 Overlap

Overlap measures the fraction of overlapping residues in a binding site for every tool (high probability residues, HPR). The values range from 0 to 1. An overlap value of 1 indicates that all residues in the binding site predicted by one tool are high probability residues that are found by the majority of the tools.

$$overlap\ score = \frac{HPR\ in\ a\ tool}{total\ HPR}$$
(5.4)

5.2.9 Correlation between Features and Quality Metrics

The correlation between features (chemical descriptors) and quality metrics was computed using the Pearson Correlation Coefficient (PCC).

$$PCC = \frac{\sum XY - \frac{(\sum X)(\sum Y)}{n}}{\sqrt{\left(\sum X^2 - \frac{(\sum X)^2}{n}\right)\left(\sum Y^2 - \frac{(\sum Y)^2}{n}\right)}}$$
(5.5)

where X is the value of a descriptor, Y is the value of the corresponding quality metrics and n is the total number of samples.

5.2.10 TScore

A tool score (TScore) was defined as the sum of the quality metrics (MCC, precision, recall, and overlap) normalized by the count of this metrics. Higher scores, in general, indicate better performance.

$$TScore = \frac{MCC + precision + recall + overlap}{4}$$
(5.6)

5.2.11 TRatio

The tool ratio (TRatio) was used to approximate the distance between a predicted TScore and the reference score. The reference score was defined separately for each protein-ligand complex as the maximum score obtained from any of the tools (maxTScore). TRatio is defined as:

$$TRatio = \frac{prediction\ score}{maxTScore}$$
(5.7)

5.2.12 Classifiers

Scikit-learn [210] was used to evaluate five classification methods Bayes [211,212], decision tree [213,214], random forest [215], support vector machine with radial kernel (svm_c_rbf), and support vector machine with linear kernel (svm_linear) [216,217]. These methods were trained on 75% of the dataset and tested on the remaining samples. As the training and test data were randomly split, we iterated the training and testing procedures 100 times to obtain more objective results.

5.3 Results

5.3.1 Pocket Distributions

The pocket tools output various quantities of pockets and formats for a given protein. DEPTH and ProACT2 output predicted pockets as a single file making it difficult to extract individual pockets for further analysis. GHECOM and PocketPicker output a predetermined number of five pockets per protein. In contrast, Fpocket, DoGSiteScorer, and IsoMif output all predicted pockets on the protein surface without limiting the number of predicted pockets for a given protein, Figure 5.3 shows the distribution of the number of predicted pockets. Fpocket, DoGSiteScorer, and IsoMif predicted 1 to 43, 3 to 112, and 13 to 147 pockets, respectively. IsoMif predicted, on average, the largest number of pockets per protein (52.1) followed by DoGSiteScorer (18.4), and Fpocket (11.7).



Figure 5.3 Distribution of the number of predicted pockets per protein for Fpocket, DoGSiteScorer, and IsoMif, respectively.

5.3.2 Size Distributions

Beside predicting different numbers of pockets, the tools also predict pockets of different sizes. For tools which provide individual pockets (Fpocket, DoGSiteScorer, IsoMif, GHECOM, and PocketPicker), we approximated the size of the predicted pockets by counting the number of residues that constitute a predicted pocket. As GHECOM and PocketPicker only output five pockets per protein, we always considered the top five pockets in each case. DoGSiteScorer, PocketPicker, and IsoMif predicted a similar number of residues for the top five pockets of 26.7, 26.5, and 23.6, respectively. Fpocket predicted 17.4 residues per pocket and GHECOM predicted the smallest number of residues per pocket, 9.0. Figure 5.4 shows the distribution of the number of residues per pocket for the tools.



Figure 5.4 Distribution of the number of pocket lining residues per pocket (top five) for Fpocket, DoGSiteScorer, IsoMif, GHECOM, and PocketPicker, respectively.

5.3.3 Prediction versus Reality

Then, we measured how well the identified pockets match the positions of bound ligands (see also discussion section). We used four quality metrics to assess the performance of each tool against a reference binding site namely Matthews correlation coefficient (MCC), precision (pre), recall (recall), and residue overlap (overlap). The Methods section provides comprehensive definitions of these metrics. Reference binding sites encompass residues within 5.0 Å distance from any atom of the small-molecule ligand bound to the protein. Figure 5.5 shows the cumulative density of each metric for each tool. In general a tool would perform well if its corresponding quality metrics are shifted to the right of the density plots.

ProACT2 and PocketPicker consistently yielded lower performance compared to the other tools. ProACT2 is on the leftmost portion of the density plots (Figure 5.5) for MCC, precision, and overlap. Similarly, PocketPicker is near or at the leftmost portion of the density plot for MCC, precision, and recall. On the other hand, GHECOM performed best (on the rightmost of the plots) in three metrics (MCC, precision, and overlap) followed by

DoGSiteScorer in recall, and DEPTH was on the rightmost of the MCC plot. On the other hand, IsoMif was placed always in the middle range.



Figure 5.5 Cumulative density of Matthews correlation coefficient (MCC), precision (pre), recall, and residues overlap (overlap).

5.3.4 Features and Quality Metrics

Next we asked whether the pocket tools work better for certain types of pockets and worse for other. Several of the tools used in this work can compute features for the predicted pockets. DoGSiteScorer, Fpocket, GHECOM, PocketPicker, and ProACT2 provide 64, 19, 5, 420, and 23 descriptors for a protein, respectively. PocketPicker [55] generates the largest number of features. However, its features focus only on shapes and buriedness of a pocket neglecting physico-chemical properties entirely. In contrast, DoGSiteScorer [56] outputs both shapes and physico-chemical descriptors. For this reason we selected DoGSiteScorer to compute features for each pocket in our reference dataset. Details regarding the features used are included in Supplementary Information Table 5.3.

Similar to the previous section, we used a cumulative density plot to visualize the correlation values. The densities are centered around zero in all metrics and tools. This suggests

predominantly weak correlations between the features and the quality metrics. Furthermore, the Pearson correlation coefficients for MCC, precision, and overlap never exceeded 0.5 in all tools. For GHECOM, the feature simple score (simpleScore) correlated slightly with recall (0.56).

5.3.5 Tool Score

To rank the tools, we defined a simple tool score (TScore) by summing MCC, precision, recall, overlap and normalized the sum by four. Figure 5.6 shows the TScores per protein-ligand complex as a heatmap. We then labelled each protein with the tool that achieved the highest TScore (maxTScore). DEPTH, GHECOM, Fpocket, DoGSiteScorer, PocketPicker, IsoMif, and ProACT2 were assigned to 65, 33, 36, 15, 0, 14, and 4 proteins, respectively. The average TScores across DEPTH (0.58), GHECOM (0.54), Fpocket (0.55), DoGSiteScorer (0.55), and IsoMif (0.53) did not differ much. Compared to this group, PocketPicker (0.40) and ProACT2 (0.36) yielded considerably lower scores. Picking the tools randomly from the seven tools yielded an average TScore of 0.53. On the other hand, if one would always pick the tool with the highest TScore for each individual protein, the maximum possible TScore is 0.67. This means that an optimal choice of tools with the highest TScore for each protein may improve the average score by 28% over choosing tools randomly.



Figure 5.6 Heatmap of TScores of DEPTH, GHECOM, Fpocket, DoGSiteScorer, PocketPicker, IsoMif, ProACT2, and maxTScore, respectively. The column maxTScore contains the maximum TScores among all methods.

5.3.6 Identifying an Optimal Pocket Tool for a Protein

Since we are now in possession of a labelled dataset along with the corresponding features, we were able to train classifiers on this dataset. The idea behind is that such classifiers could be used to identify the optimal tool for an unknown protein. We trained Bayes, decision tree, random forest, support vector machine with radial and linear kernel on 75% of the total data and tested on the remaining data. As the training and testing dataset are split randomly, we needed to account for the variance on each split, hence, we iterated the training 100 times. We found that random forest yielded an accuracy of 0.60 followed by support vector machine with radial kernel (0.55), decision tree (0.50), support vector machine with linear kernel (0.38), and Bayes (0.34).

As the TScores amongst the tools are very similar, thus, we computed the ratio of the TScores the predictions in our best performing classifier (random forest) and the maximum TScores of the reference proteins (TRatio). Such a ratio approximates the distance between the TScore of the random forest predictions and the corresponding maximum TScore of reference proteins. A ratio value near one indicates that the TScore of the random forest prediction is

very similar to the maximum TScore of the reference protein. We found that TRatios are distributed around 0.90±0.03, further indicating that the best tools have similar predictive capabilities.

As TScores are very similar across the tools we decided to output a list of tools along with the corresponding weights (probability values from the random forest model) for each protein in the prediction instead of just one optimal tool. A higher weight indicates that the classifier is confident in the decision while uniform weights across the tools indicate less confidence in the decision. For example, for EPSP synthase from *Escherichia coli* (PDB ID:2QFT) DEPTH has a clearly higher weight (0.4253) than the other tools (first row, Table 2). On the other hand, in the case of 2-naphtamidine urokinase inhibitor (PDB ID:1SQA) four tools (DEPTH, Fpocket, DoGSiteScorer, and IsoMif) have similar weights (second row, Table 2). Similarly, for the PTP16-inhibitor complex (PDB ID:2QBR) DEPTH and IsoMif have similar weights (last row, Table 2). For cases like 1SQA and 2QBR it is useful to present the predictions in the form of a list of tools instead of just one tool to allow users to interfere and make a choice based on their need or objectives. Supplementary Information Table 5.4 lists the weight for each tool obtained for each protein of the training and test sets.

Table 5.2 Examples of weight for each tool for each protein complex obtained from the classifier of the random forest model.

| PDB ID | DEPTH | GHECOM | Fpocket | DoGSiteScorer | IsoMif | ProACT2 |
|--------|--------|--------|---------|---------------|--------|---------|
| 2QFT | 0.4253 | 0.1908 | 0.1888 | 0.0744 | 0.0817 | 0.0389 |
| 1SQA | 0.2332 | 0.1385 | 0.2138 | 0.1915 | 0.1757 | 0.0472 |
| 2QBR | 0.2545 | 0.1071 | 0.1713 | 0.1459 | 0.2411 | 0.0801 |

5.4 Discussion

As mentioned in the introduction section defining the correct pockets on protein surfaces is not an easy task. This is reflected in our results. Despite similar average TScores for five tools, we found clear differences in the shape and size of the predictions (Figures 5.2 and 5.4). It was previously reported that pockets vary in shapes and can be found buried deep within a protein or narrow and shallow on the protein surface [201,202]. In addition, Villar and Kauvar [218] noted that specific amino acids such as Arg, His, Trp, and Tyr are often enriched in protein binding sites compared to the entire protein. Due to these complexities, a pocket identification tool can be very successful for a set of proteins and less successful on other sets.

On the dataset studied here and using the four quality metrics defined by us, DEPTH yielded the highest aggregated TScore for 64 proteins. However, one should note that DEPTH

predicted larger pockets compared to other tools for the majority of the proteins. This is illustrated by the number of residues per predicted pocket. DEPTH, on average, predicted pockets with 77.0 residues while other tools predicted an average smaller pockets (for top five pockets), DoGSiteScorer (26.7), PocketPicker (26.5), IsoMif (23.6), Fpocket (17.4), and GHECOM (9.0). We speculate that the larger pocket sizes might have provided an advantage to DEPTH in terms of TScores since the score is computed based on these residues.

Interestingly, when TScores were correlated to maxTScore, we found that GHECOM yielded the highest correlation (0.84) followed by DoGSiteScorer (0.78), DEPTH (0.77), and IsoMif (0.75). Not surprisingly PocketPicker (0.43) and ProACT2 (0.23) yielded the lowest correlations. Even though GHECOM did not yield the highest average TScores, the tool was either in the second or third position for each protein-ligand complex in the dataset. In addition, GHECOM also predicted the smallest pockets (Figure 5.3). The high correlation to maxTScore and the relatively small pocket size indicate that GHECOM can produce quite precise predictions and could be a reasonable choice if one does not know what tool to choose for a given protein. Alternatively, the decision for a suitable tool may of course also be based on the random forest classifier that was trained here.

5.5 Conclusions

In this work we compared seven pocket identification tools. We found that these tools predict pockets of various size and shapes. For instance, DEPTH, DoGSiteScorer and PocketPicker tend to predict larger pockets, whereas GHECOM predicts smaller sized pockets. The tools also predict various numbers of pockets per protein. IsoMif identified most pockets an average (52.1), followed by DoGSiteScorer (18.4), and Fpocket (11.7). We introduced TScores as a measure of the predictive capabilities of the tools. When one applies the optimal tool for each protein, the average TScore increases by 28% over randomly labelled proteins. The tool GHECOM constructed pockets of rather compact size and its predictions were consistently among the top three tools. Finally, we trained a random forest classifier on this dataset. The classifier outputs a list of tools with a corresponding set of weights indicating the confidence of the decision boundaries. The classifier should be useful to prioritize a set of pocket tools for protein structures that are not contained in our dataset.

Chapter 6

Conclusions and Future Works

In this thesis, we presented three different projects which mainly related to the core area of protein interactions and gene regulations. We are aimed to understanding the role of protein-protein interactions and protein-ligand interactions, protein interfaces, and pockets. We also are extending our study to gain knowledge on transcription factors targeting STIM and ORAI genes, the regulation and relation with breast cancer.

First, we performed statistical analysis on the composition of overlapping proteinprotein and protein-ligand interfaces. We started from the research question to find out whether small molecule ligands have similar physio-chemical features as protein binding interfaces when they bind at overlapping protein-protein or protein-ligand binding interfaces. We are using five different datasets from the ABC, PIBASE and TIMBAL databases. According to the statistics, we found that generally, both interface types have a very similar composition. Among the three datasets of PL complexes, we found that the protein interfaces of the Timbal dataset contain more hydrophobic residues and fewer polar residues than the two other datasets. In addition, we found that the ligands in the PIBASE and Timbal datasets bind more flat on the protein surfaces or bind deeper into pockets on the protein surface than ABC ligands. To further explore the findings, we will apply the framework on larger datasets of protein-protein and protein-ligand complexes.

For the second project, we addressed several angles about STIM and ORAI genes. We identified ten predictive interactions between transcription factors bound to promoter regions of STIM and ORAI genes based on predictions using the STRING, prePPI, and mentha databases. We used a set of genes of interest for co-expression and differential expression analysis on breast invasive carcinoma dataset which main focus on the regulation of STIM and ORAI genes. We found ORAII, ORAI2, and ORAI3 genes to be up-regulated and in contrast STIM1 and STIM2 were down-regulated. We identified several roles of STIM and ORAI genes in normal and breast cancer tissues. The results presented in this study allow us to gain knowledge and unveiling the predicted transcription factors bound to the promoter regions of STIM and ORAI genes, their regulation, and relation with breast cancer. In future, we aim to extend this study by performing gene enrichment analysis and expand this workflow to new datasets of diseases which reported related to STIM and ORAI such as Alzheimers disease and prostate cancer.

In final study presented in this thesis, we evaluate performance of seven protein pocket identification tools on a dataset of protein-ligand complexes. We analyzed how well the predicted pocket-lining residues overlap with the residues that contact the ligand. We also used the residue overlap to define a score as a measure of the predictive capabilities of the tools. We found comparable performance amongst the predictions of five tools in terms of average score. We then trained a random forest model (classifier) to output a list of suitable tools for a given protein structure to assist users in a pocket prediction scenario. This classifier should be helpful for prioritizing the tools to be used for unknown proteins or proteins that are not contained in our dataset. Future work, we will present a more comprehensive evaluation of our framework on larger datasets of protein-ligand complexes.

In summary, the studies presented in this thesis led to gain knowledge of protein interactions and regulation of STIM and ORAI genes. Overall, our future works for these three projects will involve and focusing on applying the frameworks either on new or larger datasets for better understanding and comprehensive evaluation.

REFERENCES

- Darbellay B, Arnaudeau S, König S, Jousset H, Bader C, Demaurex N, et al. STIM1and Orai1-dependent store-operated calcium entry regulates human myoblast differentiation. J Biol Chem. 2009;284: 5370–5380. doi:10.1074/jbc.M806726200
- Johnstone LS, Graham SJL, Dziadek MA. STIM proteins: Integrators of signalling pathways in development, differentiation and disease. J Cell Mol Med. 2010;14: 1890– 1903. doi:10.1111/j.1582-4934.2010.01097.x
- Flourakis M, Lehen'kyi V, Beck B, Raphaël M, Vandenberghe M, Abeele F V, et al. Orai1 contributes to the establishment of an apoptosis-resistant phenotype in prostate cancer cells. Cell Death Dis. 2010;1: e75. doi:10.1038/cddis.2010.52
- Prevarskaya N, Skryma R, Shuba Y. Calcium in tumour metastasis: new roles for known actors. Nat Rev Cancer. Nature Publishing Group; 2011;11: 609–618. doi:10.1038/nrc3105
- Gruszczynska-Biegala J, Pomorski P, Wisniewska MB, Kuznicki J. Differential roles for STIM1 and STIM2 in store-operated calcium entry in rat neurons. PLoS One. 2011;6: e19285. doi:10.1371/journal.pone.0019285
- Moccia F, Zuccolo E, Soda T, Tanzi F, Guerra G, Lisa Mapelli, et al. Stim and Orai proteins in neuronal Ca2+ signaling and excitability. Front Cell Neuroscie. 2015;9: 1–14. doi:10.3389/fncel.2015.00153
- Jemal A, Siegel R, Ward E, Murray T, Xu J, Smigal C, et al. Cancer Statistics, 2006. CA Cancer J Clin. 2006;56: 106–130. doi:10.3322/canjclin.56.2.106
- Network TCGA. Comprehensive molecular portraits of human breast tumours. Nature. 2012;490: 61–70. doi:10.1038/nature11412
- Keskin O, Tuncbag N, Gursoy A. Characterization and prediction of protein interfaces to infer protein-protein interaction networks. Curr Pharm Bi10. Jones S, Thornton JM. Principles of protein-protein interactions. Proc Natl Acad Sci USA. 1996;93: 13– 20.
- Keskin O, Nussinov R. Similar binding sites and different partners: implications to shared proteins in cellular pathways. Structure. 2007;15: 341–54. doi:10.1016/j.str.2007.01.007
- Chakrabarti P, Janin J. Dissecting protein-protein recognition sites. Proteins Struct Funct Genet. 2002;47: 334–343. doi:10.1002/prot.10085
- 13. DeLano WL. Unraveling hot spots in binding interfaces: Progress and challenges. Curr

Opin Struct Biol. 2002;12: 14–20. doi:10.1016/S0959-440X(02)00283-X

- Conte L Lo, Chothia C, Janin È. The atomic structure of protein-protein recognition sites. J Mol Biol. 1999;285: 2177–2198.
- Moreira IS, Fernandes PA, Ramos MJ. Hot spots-A review of the protein–protein interface determinant amino-acid residues. Proteins Struct Funct Bioinforma. 2007;68: 803–812. doi:10.1002/prot
- Keskin O, Ma B, Nussinov R. Hot regions in protein-protein interactions: the organization and contribution of structurally conserved hot spot residues. J Mol Biol. 2005;345: 1281–94. doi:10.1016/j.jmb.2004.10.077
- Chen Y, Kortemme T, Robertson T, Baker D, Varani G. A new hydrogen-bonding potential for the design of protein-RNA interactions predicts specific contacts and discriminates decoys. Nucleic Acids Res. 2004;32: 5147–62. doi:10.1093/nar/gkh785
- Camacho CJ, Zhang C. FastContact: rapid estimate of contact and binding free energies. Bioinformatics. 2005;21: 2534–2536. doi:10.1093/bioinformatics/bti322
- Meireles LMC, Dömling AS, Camacho CJ. ANCHOR: a web server and database for analysis of protein-protein interaction binding pockets for drug discovery. Nucleic Acids Res. 2010;38: W407-11. doi:10.1093/nar/gkq502
- Bromberg Y, Rost B. Comprehensive in silico mutagenesis highlights functionally important residues in proteins. Bioinformatics. 2008;24: i207–i212. doi:10.1093/bioinformatics/btn268
- Nooren IMA, Thornton JM. Diversity of protein-protein interactions. EMBO J. 2003;22: 3486–3492.
- 22. Acuner Ozbabacan SE, Engin HB, Gursoy A, Keskin O. Transient protein-protein interactions. Protein Eng Des Sel. 2011;24: 635–648. doi:10.1093/protein/gzr025
- Perkins JR, Diboun I, Dessailly BH, Lees JG, Orengo C. Transient Protein-Protein Interactions: Structural, Functional, and Network Properties. Structure. Elsevier Ltd; 2010;18: 1233–1243. doi:10.1016/j.str.2010.08.007
- Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, et al. The Protein Data Bank. Nucleic Acids Res. 2000;28: 235–242.
- Stark C, Breitkreutz B-J, Reguly T, Boucher L, Breitkreutz A, Tyers M. BioGRID: a general repository for interaction datasets. Nucleic Acids Res. 2006;34: D535-539. doi:10.1093/nar/gkj109
- 26. Stark C, Breitkreutz B-J, Chatr-Aryamontri A, Boucher L, Oughtred R, Livstone MS, et al. The BioGRID Interaction Database: 2011 update. Nucleic Acids Res. 2011;39:

D698-704. doi:10.1093/nar/gkq1116

- Chatr-aryamontri A, Ceol A, Palazzi LM, Nardelli G, Schneider MV, Castagnoli L, et al. MINT: the Molecular INTeraction database. Nucleic Acids Res. 2007;35: D572-4. doi:10.1093/nar/gkl950
- Bader GD, D. B, V. HCW. BIND: the Biomolecular Interaction Network Database. Nucleic Acids Res. 2003;31: 248–250. doi:10.1093/nar/gkg056
- 29. Xenarios I, Salwínski L, Duan XJ, Higney P, Kim S-M, Eisenberg D. DIP, the Database of Interacting Proteins: a research tool for studying cellular networks of protein interactions. Nucleic Acids Res. 2002;30: 303–305.
- Hermjakob H, Montecchi-Palazzi L, Lewington C, Mudali S, Kerrien S, Orchard S, et al. IntAct: an open source molecular interaction database. Nucleic Acids Res. 2004;32: D452-5. doi:10.1093/nar/gkh052
- Walter P, Ansari S, Helms V. The ABC (Analysing Biomolecular Contacts) -Database.
 J Integr Bioinform. 2007;4(1): 50.
- 32. Davis FP, Sali A. PIBASE: a comprehensive database of structurally defined protein interfaces. Bioinformatics. 2005;21: 1901–7. doi:10.1093/bioinformatics/bti277
- Higueruelo AP, Jubb H, Blundell TL. TIMBAL v2: update of a database holding small molecules modulating protein-protein interactions. Database (Oxford). 2013;2013: bat039. doi:10.1093/database/bat039
- Hu L, Benson ML, Smith RD, Lerner MG, Carlson HA. Binding MOAD (Mother Of All Databases). Proteins. 2005;60: 333–40. doi:10.1002/prot.20512
- Benson ML, Smith RD, Khazanov NA, Dimcheff B, Beaver J, Dresslar P, et al. Binding MOAD, a high-quality protein-ligand database. Nucleic Acids Res. 2008;36: D674-8. doi:10.1093/nar/gkm911
- Basse MJ, Betzi S, Bourgeas R, Bouzidi S, Chetrit B, Hamon V, et al. 2P2Idb: a structural database dedicated to orthosteric modulation of protein-protein interactions. Nucleic Acids Res. 2013;41: D824-7. doi:10.1093/nar/gks1002
- 37. Wang R, Fang X, Lu Y, Yang CY, Wang S. The PDBbind database: Methodologies and updates. J Med Chem. 2005;48: 4111–4119. doi:10.1021/jm048957q
- Liu T, Lin Y, Wen X, Jorissen RN, Gilson MK. BindingDB: a web-accessible database of experimentally determined protein-ligand binding affinities. Nucleic Acids Res. 2007;35: D198-201. doi:10.1093/nar/gkl999
- 39. Gilson MK, Liu T, Baitaluk M, Nicola G, Hwang L, Chong J. BindingDB in 2015: A public database for medicinal chemistry, computational chemistry and systems

pharmacology. Nucleic Acids Res. 2015;44: gkv1072. doi:10.1093/nar/gkv1072

- Hendlich M, Bergner A, Günther J, Klebe G. Relibase: Design and Development of a Database for Comprehensive Analysis of Protein–Ligand Interactions. J Mol Biol. 2003;326: 607–620. doi:10.1016/S0022-2836(02)01408-0
- Perot S, Sperandio O, Miteva MA, Camproux AC, Villoutreix BO. Druggable pockets and binding site centric chemical space: A paradigm shift in drug discovery. Drug Discov Today. 2010;15: 656–667. doi:10.1016/j.drudis.2010.05.015
- 42. Zheng X, Gan L, Wang E, Wang J. Pocket-based drug design: exploring pocket space. AAPS J. 2013;15: 228–41. doi:10.1208/s12248-012-9426-6
- 43. Konc J, Janezic D. Binding site comparison for function prediction and pharmaceutical discovery. Curr Opin Struct Biol. 2014;25: 34–39. doi:10.1016/j.sbi.2013.11.012
- 44. Fischer E. Einfluss der Configuration auf die Wirkung der Enzyme. Ber Dtsch Chem Ges. 1894;27: 2985–2993. doi:10.1002/cber.18940270364
- Xie Z-R, Hwang M-J. Ligand-binding site prediction using ligand-interacting and binding site-enriched protein triangles. Bioinformatics. 2012;28: 1579–85. doi:10.1093/bioinformatics/bts182
- 46. Lazo JS, Sharlow ER. Drugging undruggable molecular cancer targets. Annu Rev Pharmacol Toxicol. 2016;56: 23–40. doi:10.1146/annurev-pharmtox-010715-103440
- DeLano W. Pymol: An open-source molecular graphics tool. CCP4 Newsl Protein Crystallogr. 2002;700.
- Volkamer A, Griewel A, Grombacher T, Rarey M. Analyzing the topology of active sites: On the prediction of pockets and subpockets. J Chem Inf Model. 2010;50: 2041– 2052. doi:10.1021/ci100241y
- Kawabata T. Detection of multiscale pockets on protein surfaces using mathematical morphology. Proteins Struct Funct Bioinforma. 2010;78: 1195–1211. doi:10.1002/prot.22639
- Yu J, Zhou Y, Tanaka I, Yao M. Roll: A new algorithm for the detection of protein pockets and cavities with a rolling probe sphere. Bioinformatics. 2009;26: 46–52. doi:10.1093/bioinformatics/btp599
- Schmidtke P, Barril X. Understanding and predicting druggability. A high-throughput method for detection of drug binding sites. J Med Chem. 2010;53: 5858–67. doi:10.1021/jm100574m
- 52. Kuntz ID, Blaney JM, Oatley SJ, Langridge R, Ferrin TE. A geometric approach to macromolecule-ligand interactions. J Mol Biol. 1982;161: 269–288. doi:10.1016/0022-

2836(82)90153-X

- 53. Laskowski RA, Luscombe NM, Swindells MB, Thornton JM. Protein clefts in molecular recognition and function. Protein Sci. 1996;5: 2438–52. doi:10.1002/pro.5560051206
- Hendlich M, Rippmann F, Barnickel G. LIGSITE: Automatic and efficient detection of potential small molecule-binding sites in proteins. J Mol Graph Model. 1997;15: 359– 363. doi:10.1016/S1093-3263(98)00002-3
- 55. Weisel M, Proschak E, Schneider G. PocketPicker: analysis of ligand binding-sites with shape descriptors. Chem Cent J. 2007;1: 7. doi:10.1186/1752-153X-1-7
- Volkamer A, Kuhn D, Rippmann F, Rarey M. Dogsitescorer: A web server for automatic binding site prediction, analysis and druggability assessment. Bioinformatics. 2012;28: 2074–2075. doi:10.1093/bioinformatics/bts310
- 57. Tan KP, Varadarajan R, Madhusudhan MS. DEPTH: A web server to compute depth and predict small-molecule binding cavities in proteins. Nucleic Acids Res. 2011;39: 1–7. doi:10.1093/nar/gkr356
- 58. Tan KP, Nguyen TB, Patel S, Varadarajan R, Madhusudhan MS. Depth: a web server to compute depth, cavity sizes, detect potential small-molecule ligand-binding cavities and predict the pKa of ionizable residues in proteins. Nucleic Acids Res. 2013;41: 314–321. doi:10.1093/nar/gkt503
- Kalidas Y, Chandra N. PocketDepth: A new depth based algorithm for identification of ligand binding sites in proteins. J Struct Biol. 2008;161: 31–42. doi:10.1016/j.jsb.2007.09.005
- Laskowski RA. SURFNET: A program for visualizing molecular surfaces, cavities, and intermolecular interactions. J Mol Graph. 1995;13: 323–330. doi:10.1016/0263-7855(95)00073-9
- Kawabata T, Go N. Detection of pockets on protein surfaces using small and large probe spheres to find putative ligand binding sites. Proteins. 2007;68: 516–529. doi:10.1002/prot
- 62. Brady GP, Stouten PF. Fast prediction and visualization of protein binding pockets with PASS. J Comput Aided Mol Des. 2000;14: 383–401.
- 63. Chartier M, Najmanovich R. Detection of Binding Site Molecular Interaction Field Similarities. J Chem Inf Model. 2015;55: 1600–1615. doi:10.1021/acs.jcim.5b00333
- 64. Edelsbrunner H, Mücke EP. Three-dimensional alpha shapes. ACM Trans Graph. 1994;13: 43–72. doi:10.1145/174462.156635
- 65. Edelsbrunner H, Facello M, Liang J. On the definition and the construction of pockets

in macromolecules. Computer (Long Beach Calif). 1995; 1–16.

- Peters KP, Fauck J, Frömmel C. The automatic search for ligand binding sites in proteins of known three-dimensional structure using only geometric criteria. J Mol Biol. 1996;256: 201–213. doi:10.1006/jmbi.1996.0077
- Liang J, Edelsbrunner H, Woodward C. Anatomy of protein pockets and cavities: measurement of binding site geometry and implications for ligand design. Protein Sci. 1998;7: 1884–1897. doi:10.1002/pro.5560070905
- 68. Le Guilloux V, Schmidtke P, Tuffery P. Fpocket: an open source platform for ligand pocket detection. BMC Bioinformatics. 2009;10: 168. doi:10.1186/1471-2105-10-168
- Olsson TSG, Churchill AME, Pitt WR, Ladbury JE, Williams MA. ProACT2 Analysis of solvent accessibility, cavities and contacts in proteins and their complexes. submitted. 2012;
- Williams MA, Goodfellow JM, Thornton JM. Buried waters and internal cavities in monomeric proteins. Protein Sci. 1994;3: 1224–1235. doi:10.1002/pro.5560030808
- Laurie ATR, Jackson RM. Q-SiteFinder: An energy-based method for the prediction of protein-ligand binding sites. Bioinformatics. 2005;21: 1908–1916. doi:10.1093/bioinformatics/bti315
- 72. Glaser F, Pupko T, Paz I, Bell RE, Bechor-Shental D, Martz E, et al. Consurf: identification of functional regions in proteins by surface-mapping of phylogenetic information. Bioinformatics. 2003;19: 163–164.
- Berezin C, Glaser F, Rosenberg J, Paz I, Pupko T, Fariselli P, et al. ConSeq: The identification of functionally and structurally important residues in protein sequences. Bioinformatics. 2004;20: 1322–1324. doi:10.1093/bioinformatics/bth070
- Ashkenazy H, Erez E, Martz E, Pupko T, Ben-Tal N. ConSurf 2010: Calculating evolutionary conservation in sequence and structure of proteins and nucleic acids. Nucleic Acids Res. 2010;38: 529–533. doi:10.1093/nar/gkq399
- 75. de Rinaldis M, Ausiello G, Cesareni G, Helmer-Citterich M. Three-dimensional profiles: a new tool to identify protein surface similarities. J Mol Biol. 1998;284: 1211–21. doi:10.1006/jmbi.1998.2248
- Innis CA. SiteFiNDER|3D: A web-based tool for predicting the location of functional sites in proteins. Nucleic Acids Res. 2007;35: 489–494. doi:10.1093/nar/gkm422
- 77. Fukunishi Y, Nakamura H. Prediction of ligand-binding sites of proteins by molecular docking calculation for a random ligand library. Protein Sci. 2011;20: 95–106. doi:10.1002/pro.540

- 78. Hetenyi C, Van Der Spoel D. Blind docking of drug-sized compounds to proteins with up to a thousand residues. FEBS Lett. 2006;580: 1447–1450. doi:10.1016/j.febslet.2006.01.074
- 79. Ben-Shimon A, Niv MY. AnchorDock: Blind and flexible anchor-driven peptide docking. Structure. Elsevier Ltd; 2015;23: 929–40. doi:10.1016/j.str.2015.03.010
- Xie ZR, Hwang MJ. Ligand-binding site prediction using ligand-interacting and binding site-enriched protein triangles. Bioinformatics. 2012;28: 1579–1585. doi:10.1093/bioinformatics/bts182
- Li L, Huang Y, Xiao Y. How to use not-always-reliable binding site information in protein-protein docking prediction. PLoS One. 2013;8. doi:10.1371/journal.pone.0075936
- Del Sol Mesa A, Pazos F, Valencia A. Automatic methods for predicting functionally important residues. J Mol Biol. 2003;326: 1289–1302. doi:10.1016/S0022-2836(02)01451-1
- 83. Huang B. MetaPocket: a meta approach to improve protein ligand binding site prediction. OMICS. 2009;13: 325–30. doi:10.1089/omi.2009.0045
- Huang B, Schroeder M. LIGSITEcsc: predicting ligand binding sites using the Connolly surface and degree of conservation. BMC Struct Biol. 2006;6: 19. doi:10.1186/1472-6807-6-19
- Zhang Z, Li Y, Lin B, Schroeder M, Huang B. Identification of cavities on protein surface using multiple computational approaches for drug binding site prediction. Bioinformatics. 2011;27: 2083–2088. doi:10.1093/bioinformatics/btr331
- Levine M, Tijan R. Transcription regulation and animal diversity. Nature. 2003;424: 147–151.
- Nimwegen E van. Scaling laws in the functional content of genomes. Trends Genet.
 2003;19: 479–484. doi:10.1016/S0168-9525(03)00202-6
- Wasserman WW, Sandelin A. Applied bioinformatics for the identification of regulatory elements. Nat Rev Genet. 2004;5: 276–287. doi:10.1038/nrg1315
- Dreos R, Ambrosini G, Cavin Perier R, Bucher P. EPD and EPDnew, high-quality promoter resources in the next-generation sequencing era. Nucleic Acids Res. 2013;41: D157–D164. doi:10.1093/nar/gks1233
- 90. Bryne JC, Valen E, Tang M-HE, Marstrand T, Winther O, da Piedade I, et al. JASPAR, the open access database of transcription factor-binding profiles: new content and tools in the 2008 update. Nucleic Acids Res. 2008;36: D102-6. doi:10.1093/nar/gkm955

- 91. Kulakovskiy I V, Medvedeva YA, Schaefer U, Kasianov AS, Vorontsov IE, Bajic VB, et al. HOCOMOCO: a comprehensive collection of human transcription factor binding sites models. Nucleic Acids Res. 2013;41: D195-202. doi:10.1093/nar/gks1089
- 92. Matys V, Fricke E, Geffers R, Gößling E, Haubrock M, Hehl R, et al. TRANSFAC®: Transcriptional regulation, from patterns to profiles. Nucleic Acids Res. 2003;31: 374– 378. doi:10.1093/nar/gkg108
- 93. Kilpatrick AM. A novel stochastic and entropy-based Expection-Maximisation algorithm for transcription factor binding site motif discovery. 2014. doi:10.1007/s13398-014-0173-7.2
- 94. D'haeseleer P. What are DNA sequence motifs? Nat Biotechnol. 2006;24: 423–425. doi:10.1038/nbt0406-423
- 95. Rapaport F, Khanin R, Liang Y, Pirun M, Krek A, Zumbo P, et al. Comprehensive evaluation of differential gene expression analysis methods for RNA-seq data. Genome Biol. 2013;14: R95. doi:10.1186/gb-2013-14-9-r95
- 96. Robinson MD, McCarthy DJ, Smyth GK. edgeR: A Bioconductor package for differential expression analysis of digital gene expression data. Bioinformatics. 2009;26: 139–140. doi:10.1093/bioinformatics/btp616
- 97. Anders S, Huber W, Nagalakshmi U, Wang Z, Waern K, Shou C, et al. Differential expression analysis for sequence count data. Genome Biol. 2010;11: R106. doi:10.1186/gb-2010-11-10-r106
- 98. Varet H, Coppée J-Y, Dillies M-A. SARTools: a DESeq2- and edgeR-based R pipeline for comprehensive differential analysis of RNA-Seq data. PLoS One. 2016;11. doi:10.1371/journal.pone.0157022
- 99. Anders S, McCarthy DJ, Chen Y, Okoniewski M, Smyth GK, Huber W, et al. Countbased differential expression analysis of RNA sequencing data using R and Bioconductor. Nat Protoc. 2013;8: 1765–1786. doi:10.1038/nprot.2013.099
- 100. Dillies MA, Rau A, Aubert J, Hennequet-Antier C, Jeanmougin M, Servant N, et al. A comprehensive evaluation of normalization methods for Illumina high-throughput RNA sequencing data analysis. Brief Bioinform. 2013;14: 671–683. doi:10.1093/bib/bbs046
- 101. Benjamin Y, Hochberg Y. Controling the False Discovery Rate: A Practical and POwerful Approach to Multiple Testing. J R Stat Soc Ser C (Applied Stat. 1995;57: 289–300. doi:10.1038/203024b0
- 102. Yu H, Liu B-H, Ye Z-Q, Li C, Li Y-X, Li Y-Y. Link-based quantitative methods to identify differentially coexpressed genes and gene pairs. BMC Bioinformatics. 2011;12:

315. doi:10.1186/1471-2105-12-315

- 103. Amar D, Safer H, Shamir R. Dissection of regulatory networks that are altered in disease via differential co-expression. PLoS Comput Biol. 2013;9. doi:10.1371/journal.pcbi.1002955
- 104. Tesson BM, Breitling R, Jansen RC. DiffCoEx: a simple and sensitive method to find differentially coexpressed gene modules. BMC Bioinformatics. 2010;11: 497. doi:10.1186/1471-2105-11-497
- 105. Reverter A, Ingham A, Lehnert SA, Tan SH, Wang Y, Ratnakumar A, et al. Simultaneous identification of differential gene expression and connectivity in inflammation, adipogenesis and cancer. Bioinformatics. 2006;22: 2396–2404. doi:10.1093/bioinformatics/btl392
- 106. Langfelder P, Horvath S. WGCNA: an R package for weighted correlation network analysis. BMC Bioinformatics. 2008;9: 559. doi:10.1186/1471-2105-9-559
- 107. DiLeo M V., Strahan GD, den Bakker M, Hoekenga OA. Weighted correlation network analysis (WGCNA) applied to the tomato fruit metabolome. PLoS One. 2011;6. doi:10.1371/journal.pone.0026683
- 108. Mohamed R, Degac J, Helms V. Composition of overlapping protein-protein and protein-ligand interfaces. PLoS One. 2015;10: 4–6. doi:10.1371/journal.pone.0140965
- 109. Nooren IMA, Thornton JM. Structural characterisation and functional significance of transient protein–protein Interactions. J Mol Biol. 2003;325: 991–1018. doi:10.1016/S0022-2836(02)01281-0
- 110. Jones S, Daley DT, Luscombe NM, Berman HM, Thornton JM. Protein-RNA interactions: a structural analysis. Nucleic Acids Res. 2001;29: 943–54.
- 111. Petrescu A-J, Milac A-L, Petrescu SM, Dwek RA, Wormald MR. Statistical analysis of the protein environment of N-glycosylation sites: implications for occupancy, structure, and folding. Glycobiology. 2004;14: 103–14. doi:10.1093/glycob/cwh008
- 112. Lejeune D, Delsaux N, Charloteaux B, Thomas A, Brasseur R. Protein-nucleic acid recognition: statistical analysis of atomic interactions and influence of DNA structure. Proteins. 2005;61: 258–71. doi:10.1002/prot.20607
- 113. Davis FP, Sali A. The Overlap of Small Molecule and Protein Binding Sites within Families of Protein Structures. PLoS Comput Biol. 2010;6. doi:10.1371/journal.pcbi.1000668
- 114.Walter P, Metzger J, Thiel C, Helms V. Predicting where small molecules bind at
protein-protein interfaces.PLoSOne.2013;8:e58583.

doi:10.1371/journal.pone.0058583

- 115. Koes DR, Camacho CJ. Small-molecule inhibitor starting points learned from proteinprotein interaction inhibitor structure. Bioinformatics. 2012;28: 784–91. doi:10.1093/bioinformatics/btr717
- Huang Y, Niu B, Gao Y, Fu L, Li W. CD-HIT Suite: a web server for clustering and comparing biological sequences. Bioinformatics. 2010;26: 680–2. doi:10.1093/bioinformatics/btq003
- 117. Li W, Godzik A. Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. Bioinformatics. 2006;22: 1658–9. doi:10.1093/bioinformatics/btl158
- Holland RCG, Down TA, Pocock M, Prlić A, Huen D, James K, et al. BioJava: an opensource framework for bioinformatics. Bioinformatics. 2008;24: 2096–7. doi:10.1093/bioinformatics/btn397
- 119. Hubbard S, Thornton JM. Naccess. v2.1.1 [Internet]. Manchester U.K; 1996.
- Eisenberg D, Weiss RM, Terwilliger TC, Wilcox W. Hydrophobic moments and protein structure. Faraday Symp Chem Soc. 1982;17: 109. doi:10.1039/fs9821700109
- 121. Higueruelo AP, Schreyer A, Bickerton GRJ, Blundell TL, Pitt WR. What can we learn from the evolution of protein-ligand interactions to aid the design of new therapeutics? PLoS One. 2012;7: e51742. doi:10.1371/journal.pone.0051742
- Eyrisch S, Helms V. Transient pockets on protein surfaces involved in protein-protein interaction. J Med Chem. 2007;50: 3457–64. doi:10.1021/jm070095g
- Yan C, Wu F, Jernigan RL, Dobbs D, Hanovar V. Characterization of Protein-protein Interfaces. Protein J. 2008;27: 59–70. doi:10.1007/s10930-007-9108-x.Characterization
- 124. Bickerton GR, Higueruelo AP, Blundell TL. Comprehensive, atomic-level characterization of structurally characterized protein-protein interactions: the PICCOLO database. BMC Bioinformatics. BioMed Central Ltd; 2011;12: 313. doi:10.1186/1471-2105-12-313
- 125. Janin J, Bahadur RP, Chakrabarti P. Protein-protein interaction and quaternary structure.
 Q Rev Biophys. 2008;41: 133–180. doi:10.1017/S0033583508004708
- 126. Glaser F, Steinberg DM, Vakser IA, Ben-tal N. Residue Frequencies and Pairing Preferences at Protein-protein interfaces. PROTEINS Struct Funct Genet. 2001;43: 89– 102.
- 127. Jones S, Thornton JM. Analysis of protein-protein interaction sites using surface patches. J Mol Biol. 1997;272: 121–32. doi:10.1006/jmbi.1997.1234

- 128. Jin L, Wang W, Fang G. Targeting protein-protein interaction by small molecules. Annu Rev Pharmacol Toxicol. 2014;54: 435–56. doi:10.1146/annurev-pharmtox-011613-140028
- 129. Bogan AA, Thorn KS. Anatomy of hot spots in protein interfaces. J Mol Biol. 1998;280:
 1–9. doi:10.1006/jmbi.1998.1843
- Jones S, Marin A, Thornton JM. Protein domain interfaces: characterization and comparison with oligomeric protein interfaces. Protein Eng. 2000;13: 77–82. doi:10.1093/protein/13.2.77
- Talavera D, Robertson DL, Lovell SC. Characterization of protein-protein interaction interfaces from a single species. PLoS One. 2011;6: e21053. doi:10.1371/journal.pone.0021053
- 132. McCoy AJ, Chandana Epa V, Colman PM. Electrostatic complementarity at protein/protein interfaces. J Mol Biol. 1997;268: 570–584. doi:10.1006/jmbi.1997.0987
- Sheinerman FB, Norel R, Honig B. Electrostatic aspects of protein-protein interactions. Curr Opin Struct Biol. 2000;10: 153–9.
- 134. Ofran Y, Rost B. Analysing six types of protein-protein interfaces. J Mol Biol. 2003;325:
 377–387. doi:10.1016/S0022-2836(02)01223-8
- 135. Chen K, Kurgan L. Investigation of atomic level patterns in protein--small ligand interactions. PLoS One. 2009;4: e4473. doi:10.1371/journal.pone.0004473
- 136. Hakulinen R, Puranen S, Lehtonen J V, Johnson MS, Corander J. Probabilistic prediction of contacts in protein-ligand complexes. PLoS One. 2012;7: e49216. doi:10.1371/journal.pone.0049216
- 137. Bondi A. van der Waals Volumes and Radii. J Phys Chem. 1965;68: 441-451.
- 138. Parekh AB, Penner R. Store depletion and calcium influx. Physiol Rev. 1997;77: 901– 30.
- 139. Putney JW. A model for receptor-regulated Calcium entry. Cell Calcium. 1986;7: 1–12.
- 140. Liou J, Kim ML, Won DH, Jones JT, Myers JW, Ferrell JE, et al. STIM is a Ca2+ sensor essential for Ca2+-store- depletion-triggered Ca2+ influx. Curr Biol. 2005;15: 1235–1241. doi:10.1016/j.cub.2005.05.055
- 141. Roos J, DiGregorio PJ, Yeromin A V., Ohlsen K, Lioudyno M, Zhang S, et al. STIM1, an essential and conserved component of store-operated Ca 2+ channel function. J Cell Biol. 2005;169: 435–445. doi:10.1083/jcb.200502019
- 142. Zhang SL, Yu Y, Roos J, Kozak JA, Deerinck TJ, Ellisman MH, et al. STIM1 is a Ca2+ sensor that activates CRAC channels and migrates from the Ca2+ store to the plasma

membrane. Nature. 2005;437: 902-5. doi:10.1038/nature04147

- 143. Prakriya M, Feske S, Gwack Y, Srikanth S, Rao A, Hogan PG. Orai1 is an essential pore subunit of the CRAC channel. Nature. 2006;443: 230–233. doi:10.1038/nature05122
- 144. Yeromin A V, Zhang SL, Jiang W, Yu Y, Safrina O, Cahalan MD. Molecular identification of the CRAC channel by altered ion selectivity in a mutant of Orai. Nature. 2006;443: 226–9. doi:10.1038/nature05108
- 145. Lis A, Peinelt C, Beck A, Parvez S, Monteilh-Zoller M, Fleig A, et al. CRACM1, CRACM2, and CRACM3 are store-operated Ca2+ channels with distinct functional properties. Curr Biol. 2007;17: 794–800. doi:10.1016/j.cub.2007.03.065
- 146. Mercer JC, DeHaven WI, Smyth JT, Wedel B, Boyles RR, Bird GS, et al. Large storeoperated calcium selective currents due to co-expression of Orai1 or Orai2 with the intracellular calcium sensor, Stim1. J Biol Chem. 2006;281: 24979–24990. doi:10.1074/jbc.M604589200
- 147. Rosado JA, Diez R, Smani T, Jardín I. STIM and Orai1 variants in store-operated calcium entry. Front Pharmacol. 2016;6: 1–9. doi:10.3389/fphar.2015.00325
- 148. Williams RT, Manji SSM, Parker NJ, Hancock MS, Stekelenburg L Van, Eid J, et al. Identification and characterization of the STIM (stromal interaction molecule) gene family : coding for a novel class of transmembrane proteins. Biochem J. 2001;357: 673– 685.
- Dziadek MA, Johnstone LS. Biochemical properties and cellular localisation of STIM proteins. Cell Calcium. 2007;42: 123–132. doi:10.1016/j.ceca.2007.02.006
- Yuan JP, Zeng W, Dorwart MR, Choi Y-J, Worley PF, Muallem S. SOAR and the polybasic STIM1 domains gate and regulate Orai channels. Nat Cell Biol. 2009;11: 337– 43. doi:10.1038/ncb1842
- 151. Thiel M, Lis A, Penner R. STIM2 drives Ca2+ oscillations through store-operated Ca2+ entry caused by mild store depletion. J Physiol. 2013;591: 1433–45. doi:10.1113/jphysiol.2012.245399
- Feske S, Skolnik EY, Prakriya M. Ion channels and transporters in lymphocyte function and immunity. Nat Rev Immunol. Nature Publishing Group; 2012;12: 532–47. doi:10.1038/nri3233
- Hoth M, Niemeyer BA. The neglected CRAC proteins: Orai2, Orai3, and STIM2. Curr Top Membr. 2013;71: 237–271.
- 154. Vig M, Peinelt C, Beck A, Koomoa DL, Rabah D, Koblan-Huberson M, et al. CRACM1Is a Plasma Membrane Protein Essential for Store-Operated Ca2+ Entry. Science (80-).

2006;312: 1120–1223. doi:10.1007/s13398-014-0173-7.2

- 155. Derler I, Jardin I, Romanin C. The molecular mechanisms of STIM/Orai communication. Am J Physiol Cell Physiol. 2016;310: C643–C662. doi:10.1152/ajpcell.00007.2016
- 156. Muik M, Fahrner M, Schindl R, Stathopulos P, Frischauf I, Derler I, et al. STIM1 couples to ORAI1 via an intramolecular transition into an extended conformation. EMBO J. Nature Publishing Group; 2011;30: 1678–1689. doi:10.1038/emboj.2011.79
- 157. Muik M, Frischauf I, Derler I, Fahrner M, Bergsmann J, Eder P, et al. Dynamic coupling of the putative coiled-coil domain of ORAI1 with STIM1 mediates ORAI1 channel activation. J Biol Chem. 2008;283: 8014–8022. doi:10.1074/jbc.M708898200
- 158. Berna-Erro A, Woodard GE, Rosado JA. Orais and STIMs: Physiological mechanisms and disease. J Cell Mol Med. 2012;16: 407–424. doi:10.1111/j.1582-4934.2011.01395.x
- 159. Bojarski L, Pomorski P, Szybinska A, Drab M, Skibinska-Kijek A, Gruszczynska-Biegala J, et al. Presenilin-dependent expression of STIM proteins and dysregulation of capacitative Ca2+ entry in familial Alzheimer's disease. Biochim Biophys Acta - Mol Cell Res. 2009;1793: 1050–1057. doi:10.1016/j.bbamcr.2008.11.008
- 160. Schuhmann MK, Stegner D, Berna-Erro A, Bittner S, Braun A, Kleinschnitz C, et al. Stromal interaction molecules 1 and 2 are key regulators of autoreactive T cell activation in murine autoimmune central nervous system inflammation. J Immunol. 2010;184: 1536–1542. doi:10.4049/jimmunol.0902161
- 161. Cheng KT, Alevizos I, Liu X, Swaim WD, Yin H, Feske S, et al. STIM1 and STIM2 protein deficiency in T lymphocytes underlies development of the exocrine gland autoimmune disease, Sjogren's syndrome. Proc Natl Acad Sci. 2012;109: 14544–14549. doi:10.1073/pnas.1207354109
- 162. Aytes A, Mollevi DG, Martinez-Iniesta M, Nadal M, Vidal A, Morales A, et al. Stromal interaction molecule 2 (STIM2) is frequently overexpressed in colorectal tumors and confers a tumor cell growth suppressor phenotype. Mol Carcinog. 2012;51: 746–753. doi:10.1002/mc.20843
- 163. Stanisz H, Saul S, Müller CSL, Kappl R, Niemeyer B a., Vogt T, et al. Inverse regulation of melanoma growth and migration by Orai1/STIM2-dependent calcium entry. Pigment Cell Melanoma Res. 2014;27: 442–453. doi:10.1111/pcmr.12222
- 164. Sun S, Zhang H, Liu J, Popugaeva E, Xu NJ, Feske S, et al. Reduced synaptic STIM2 expression and impaired store-operated calcium entry cause destabilization of mature spines in mutant presenilin mice. Neuron. Elsevier Inc.; 2014;82: 79–93.
doi:10.1016/j.neuron.2014.02.019

- 165. Yang S, Zhang JJ, Huang XY. Orai1 and STIM1 are critical for breast tumor cell migration and metastasis. Cancer Cell. Elsevier Inc.; 2009;15: 124–134. doi:10.1016/j.ccr.2008.12.019
- 166. Motiani RK, Abdullaev IF, Trebak M. A novel native store-operated calcium channel encoded by Orai3: selective requirement of Orai3 versus Orai1 in estrogen receptorpositive versus estrogen receptor-negative breast cancer cells. J Biol Chem. 2010;285: 19173–19183. doi:10.1074/jbc.M110.102582
- 167. McAndrew D, Grice DM, Peters A a, Smart CE, Brown M a, Kenny P a, et al. ORAI1-Mediated Calcium Influx in Lactation and in Breast Cancer ORAI1-Mediated Calcium Influx in Lactation and in Breast Cancer. Mol Cancer Ther. 2011;10: 448–460. doi:10.1158/1535-7163.MCT-10-0923
- 168. Faouzi M, Hague F, Potier M, Ahidouch A, Sevestre H, Ouadid-Ahidouch H. Downregulation of Orai3 arrests cell-cycle progression and induces apoptosis in breast cancer cells but not in normal breast epithelial cells. J Cell Physiol. 2011;226: 542–551. doi:10.1002/jcp.22363
- Bhardwaj R, Hediger MA, Demaurex N. Redox modulation of STIM-ORAI signaling. Cell Calcium. Elsevier Ltd; 2016; doi:10.1016/j.ceca.2016.03.006
- Hannenhalli S. Eukaryotic transcription factor binding sites Modeling and integrative search methods. Bioinformatics. 2008;24: 1325–1331. doi:10.1093/bioinformatics/btn198
- 171. Lachmann A, Xu H, Krishnan J, Berger SI, Mazloom AR, Maayan A. ChEA: Transcription factor regulation inferred from integrating genome-wide ChIP-X experiments. Bioinformatics. 2010;26: 2438–2444. doi:10.1093/bioinformatics/btq466
- 172. Wang J, Zhuang J, Iyer S, Lin XY, Greven MC, Kim BH, et al. Factorbook.org: A Wikibased database for transcription factor-binding data generated by the ENCODE consortium. Nucleic Acids Res. 2013;41: 171–176. doi:10.1093/nar/gks1221
- 173. Feingold E, Good P, Guyer M, Kamholz S, Liefer L, Wetterstrand K, et al. The ENCODE (ENCyclopedia of DNA elements) Project. 2004;306: 636–640. doi:10.1126/science.1105136
- 174. Sandelin A, Alkema W, Engström P, Wasserman WW, Lenhard B. JASPAR: an openaccess database for eukaryotic transcription factor binding profiles. Nucleic Acids Res. 2004;32: 91D–94. doi:10.1093/nar/gkh012
- 175. Xia K, Dong D, Han J-DJ. IntNetDB v1.0: an integrated protein-protein interaction

network database generated by a probabilistic model. BMC Bioinformatics. 2006;7: 508. doi:10.1186/1471-2105-7-508

- 176. Calderone A, Castagnoli L, Cesareni G. Mentha: a Resource for Browsing Integrated Protein-Interaction Networks. Nat Methods. Nature Publishing Group; 2013;10: 690. doi:10.1038/nmeth.2561
- 177. Brown KR, Jurisica I. Online predicted human interaction database. Bioinformatics. 2005;21: 2076–2082. doi:10.1093/bioinformatics/bti273
- 178. Zhang QC, Petrey D, Garzón JI, Deng L, Honig B. PrePPI: A structure-informed database of protein-protein interactions. Nucleic Acids Res. 2013;41: 828–833. doi:10.1093/nar/gks1231
- 179. Mellor JC, Yanai I, Clodfelter KH, Mintseris J, DeLisi C. Predictome: a database of putative functional links between proteins. Nucleic Acids Res. 2002;30: 306–9. doi:10.1093/nar/30.1.306
- McDowall MD, Scott MS, Barton GJ. PIPs: human protein-protein interaction prediction database. Nucleic Acids Res. 2009;37: D651-6. doi:10.1093/nar/gkn870
- 181. Franceschini A, Szklarczyk D, Frankild S, Kuhn M, Simonovic M, Roth A, et al. STRING v9.1: protein-protein interaction networks, with increased coverage and integration. Nucleic Acids Res. 2013;41: D808-15. doi:10.1093/nar/gks1094
- 182. Szklarczyk D, Franceschini A, Kuhn M, Simonovic M, Roth A, Minguez P, et al. The STRING database in 2011: functional interaction networks of proteins, globally integrated and scored. Nucleic Acids Res. 2011;39: D561-8. doi:10.1093/nar/gkq973
- 183. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, et al. Cytoscape : A Software Environment for Integrated Models of Biomolecular Interaction Networks. Genome Res. 2003;13: 2498–2504. doi:10.1101/gr.1239303
- Périer RC, Junier T, Bonnard C, Bucher P. The eukaryotic promoter database (EPD). Nucleic Acids Res. 1999;27: 307–309. doi:10.1093/nar/27.1.307
- Grant CE, Bailey TL, Noble WS. FIMO: scanning for occurrences of a given motif. Bioinformatics. 2011;27: 1017–1018. doi:10.1093/bioinformatics/btr064
- Bailey TL, Johnson J, Grant CE, Noble WS. The MEME Suite. Nucleic Acids Res. 2015;43: W39–W49. doi:10.1093/nar/gkv416
- 187. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. Genome Biol. 2014;15: 550. doi:10.1186/s13059-014-0550-8
- 188. Eylenstein A, Schmidt S, Gu S, Yang W, Schmid E, Schmidt EM, et al. Transcription

factor NF-κB regulates expression of pore-forming Ca 2+ channel unit, Orai1, and its activator, STIM1, to control Ca 2+ entry and affect cellular functions. J Biol Chem. 2012;287: 2719–2730. doi:10.1074/jbc.M111.275925

- 189. Chang W-C, Fang Y-Y, Chang H-W, Chuang L-Y, Lin Y-D, Hou M-F, et al. Identifying association model for single-nucleotide polymorphisms of ORAI1 gene for breast cancer. Cancer Cell Int. Cancer Cell International; 2014;14: 29. doi:10.1186/1475-2867-14-29
- 190. Hoth M. CRAC channels, calcium, and cancer in light of the driver and passenger concept. Biochim Biophys Acta - Mol Cell Res. Elsevier B.V.; 2016;1863: 1408–1417. doi:10.1016/j.bbamcr.2015.12.009
- Dimova DK, Dyson NJ. The E2F transcriptional network: old acquaintances with new faces. Oncogene. 2005;24: 2810–2826. doi:10.1038/sj.onc.1208612
- 192. DeGregori J, Johnson DG. Distinct and overlapping roles for E2F family members in transcription, proliferation and apoptosis. Curr Mol Med. 2006;6: 739–748. doi:10.1038/nature08544
- 193. Oeckinghaus A, Ghosh S. The NF-kappaB family of transcription factors and its regulation. Cold Spring Harb Perspect Biol. 2009;1: 1–14. doi:10.1101/cshperspect.a000034
- 194. Dragoni S, Laforenza U, Bonetti E, Lodola F, Bottino C, Berra-Romani R, et al. Vascular endothelial growth factor stimulates endothelial colony forming cells proliferation and tubulogenesis by inducing oscillations in intracellular Ca2+ concentration. Stem Cells. 2011;29: 1898–907. doi:10.1002/stem.734
- 195. Parekh AB. Decoding cytosolic Ca2+ oscillations. Trends Biochem Sci. Elsevier Ltd;
 2011;36: 78–87. doi:10.1016/j.tibs.2010.07.013
- 196. Clapham DE. Calcium Signaling. Cell. 2007;131: 1047–1058. doi:10.1016/j.cell.2007.11.028
- 197. Levy C, Khaled M, Fisher DE. MITF: master regulator of melanocyte development and melanoma oncogene. Trends Mol Med. 2006;12: 406–414. doi:10.1016/j.molmed.2006.07.008
- 198. Stanisz H, Vultur A, Herlyn M, Roesch A, Bogeski I. The role of Orai / STIM calcium channels in melanocytes and melanoma. J Physiol. 2016;0: 1–11. doi:10.1113/JP271141.This
- 199. Sotriffer C, Klebe G. Identification and mapping of small-molecule binding sites in proteins: computational tools for structure-based drug design. II Farm. 2002;57: 243–

251. doi:10.1016/S0014-827X(02)01211-9

- 200. Henrich S, Salo-Ahen OMH, Huang B, Rippmann FF, Cruciani G, Wade RC. Computational approaches to identifying and characterizing protein binding sites for ligand design. J Mol Recognit. 2009;23: 209–19. doi:10.1002/jmr.984
- Coleman RG, Sharp KA. Protein pockets: inventory, shape, and comparison. J Chem Inf Model. 2010;50: 589–603. doi:10.1021/ci900397t
- 202. Levitt M, Park BH. Water: now you see it, now you don't. Curr Biol Struct. 1993;1: 223–226.
- 203. Pupko T, Bell RE, Mayrose I, Glaser F, Ben-Tal N. Rate4Site: an algorithmic tool for the identification of functional regions in proteins by surface mapping of evolutionary determinants within their homologues. Bioinformatics. 2002;18 Suppl 1: S71–S77. doi:10.1093/bioinformatics/18.suppl_1.S71
- 204. Oliveira SHP, Ferraz FA, Honorato R V, Xavier-Neto J, Sobreira TJP, de Oliveira PSL. KVFinder: steered identification of protein cavities as a PyMOL plugin. BMC Bioinformatics. 2014;15: 197. doi:10.1186/1471-2105-15-197
- 205. Degac J, Winter U, Helms V. Graph-based clustering of predicted ligand-binding pockets on protein surfaces. J Chem Inf Model. 2015;55: 1944–1952. doi:10.1021/acs.jcim.5b00045
- 206. Wang R, Fang X, Lu Y, Wang S. The PDBbind database: collection of binding affinities for protein-ligand complexes with known three-dimensional structures. J Med Chem. 2004;47: 2977–80. doi:10.1021/jm0305801
- 207. Matthews BW. Comparison of the predicted and observed secondary structure of T4 phage lysozyme. BBA Protein Struct. 1975;405: 442–451. doi:10.1016/0005-2795(75)90109-9
- 208. Murakami Y, Mizuguchi K. Applying the Naive Bayes classifier with kernel density estimation to the prediction of protein-protein interaction sites. Bioinformatics. 2010;26: 1841–1848. doi:10.1093/bioinformatics/btq302
- 209. Baldi P, Brunak S, Chauvin Y, Andersen C. AF, Nielsen H. Assessing the accuracy of prediction algorithms for classification: an overview. Bioinformatics. 2000;16: 412– 424. doi:10.1093/bioinformatics/16.5.412
- 210. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikitlearn: machine learning in Python. J Mach Learn Res. 2011;12: 2825–2830. doi:10.1007/s13398-014-0173-7.2
- 211. Abdo A, Chen B, Mueller C, Salim N, Willett P. Ligand-based virtual screening using

Bayesian inference networks. J Chem Inf Model. 2010;50: 1012–1020.

- 212. Ahmed A, Abdo A, Salim N. Ligand-based virtual screening using Bayesian inference network and reweighted fragments. Sci World J. 2012;2012: 1–7. doi:10.1100/2012/410914
- 213. Chen X, Wang M, Zhang H. The use of classification trees for bioinformatics. Wiley Interdiscip Rev Data Min Knowl Discov. 2011;1: 55–63. doi:10.1002/widm.14
- Zahiri J, Bozorgmehr JH, Masoudi-Nejad A. Computational prediction of proteinprotein interaction networks: algorithms and resources. Curr Genomics. 2013;14: 397– 414. doi:10.2174/1389202911314060004
- 215. Breiman L. Random forests. Mach Learn. 2001;45: 5–32. doi:10.1023/A:1010933404324
- 216. Cortes C, Vapnik V. Support-vector networks. Mach Learn. 1995;20: 273–297. doi:10.1007/BF00994018
- Jorissen RN, Gilson MK. Virtual screening of molecular databases using a Support Vector Machine. J Chem Inf Model. 2005;45: 549–561.
- 218. Villar HO, Kauvar LM. Amino acid preferences at protein binding sites. FEBS Lett. 1994;349: 125–30. doi:10.1016/0014-5793(94)00648-2

Supplementary Information

8.1 Supplementary Information for Chapter 3

Table 3.1 Datasets of PP and PL complexes.

Dataset of PP complexes from the ABC database (Table A).

PDB ID Chain Chain 15C8 L Н 1HE8 В A 1R5T А В 2BL0 В А В Е Т А В A 2BTO A 1A09 1HFY А 1R8Q С 1A6U Н L 1HJA В L 1RFX В 2C3N А В В D С В С 1AB8 А 1HQK Е 1RH7 2CHP D F А А В А А В 1AHW 1HWU В 1RQ7 2COG 1ASL В А 1130 В Е 1S4Y В А 2CUY В А С В A 1AV1 В А 118F 1S7Y В 2D26 А В 1AVA А С 1ICF А L 1SC1 А В 2D4V В А В G A В А G В 1AY7 1IFV 1SGF 2DSQ L А С В Е С В 1AZZ 1IIN А 1SGR Т 2DVG A В С В В A В 2EV4 A 1B06 1JSU 1SND F С 1JTH Е 1B99 D А 1STF L 2FHZ В A С С F С 1BJQ A 1JZD В 1SUV D 2FR6 D F В В 1BMF 1KQD А 1T6G С А 2G2U А В С С Κ Н 1BRR А 1KQM А 1TH7 L 2GD4 L А С В A А В L 1BVI 1KWS 1U3R 2GMR Μ Ρ Т 1BVN 1MCI В А 1U3W В A 2GPV D В 1BZX Е I 1MCV А L 1UGH Е L 2GQD В А D A В A A В В A 1C2O 1MI3 1UVC 2GVM Е F В A В С В A 1CKG 1MSA 1VGQ 2H1L 1CYY А В 1N8O С Е 1VLZ В А 2HI7 А В D В С В A Н D 1D7F A 1N9S 1W0I 212R 1DFJ Е L 1NHG В D 1W29 С D 2IDO А В 1DLE А В 1NMA Ν Η 1X1U А D 2IJO А Т В D 1NW9 А В 1X1Z А В 2IPJ В A 1DM5 В С В А Е 1DPJ A 1NYS А 1XBY 2IWG D Е J А В В А В A 1E0F 2IYG 106S 1XMZ 1E8U В А 1081 В А 1XX9 В D 2J3K А В Н А 1Y48 В А 1EYS Μ 10BQ В Е L 2J96 1EZL С D 10KT В А 1YGP А В 2NX5 А D А В В С Н G Н 1F45 10ME А 1YI5 2NY7 в В 1F46 А 1009 А 1YRQ А Н 208A А Т 1F88 В А 10PF В А 1ZJD А В 20CC А В Е С 1FAK Н L 1P0S Н 1ZRS В А 20L1 А С 1FI8 А 1P4I L Н 1ZUX D В 2P1L А В С 1FX9 В A 1P7L А В 2AAX А В 2PRG А В 1G0A D 1PA0 В А 2AGY В D 2PVO D А 1G21 Н G 1PMO D С 2AMT С В 2TNF В А 1G3N А В 1PPF Е L 2AQ3 G А 3PCD А Μ Κ Н 1G6V А 1PST Μ 2AY5 А В 4AKE А В 1GIF В А 1PYG В А В А 6PFK С D 2AYQ 1GQ3 С В 1QGE D Е В А 7TIM В А 2B42 1GXS D С 1QPX А В 2BE6 А D Е А 1H6A В А 1QRN D 2BEX С

Dataset of PL complexes from the ABC database (Table B).

| PDB ID | Chain | Ligand | | | | | | | | | |
|--------|-------|--------|------|---|-----|------|---|-----|------|---|-----|
| 1A05 | В | IPM | 1IT6 | А | CYU | 1TI1 | А | D12 | 2DQV | А | GAL |
| 1A0J | А | BEN | 1IZ2 | А | SUM | 1TR5 | А | THP | 2FMH | А | TRS |
| 1A8J | L | PME | 1JI5 | В | MPD | 1TXC | А | 2AN | 2FNW | А | REP |
| 1ANK | А | AMP | 1JTK | А | THU | 1U0H | В | ONM | 2FYD | А | PG4 |
| 1AR1 | А | LDA | 1KJ1 | D | MAN | 1U3T | В | ССВ | 2G2Z | А | COZ |
| 1BG9 | А | BGC | 1KMH | В | TTX | 1UBH | S | MPD | 2G7Y | А | MO9 |
| 1BIW | В | S80 | 1KYN | А | KTP | 1USR | А | SIA | 2GJ6 | D | 3IB |
| 1BLC | А | CEM | 1L7Z | А | MYR | 1UTM | А | PEA | 2GOO | С | NDG |
| 1BMQ | А | MNO | 1L9B | L | HTO | 1UX0 | А | THU | 2H0T | А | EPE |
| 1BQI | А | SBA | 1L9H | А | HTO | 1V3V | В | 50P | 2H6Y | А | MPD |
| 1BWO | А | LPC | 1LBC | А | CYZ | 1V84 | А | NAG | 2HA3 | А | P6G |
| 1C50 | А | CHI | 1LIN | А | TFP | 1W5F | В | G2P | 2HG8 | А | MLE |
| 1CGY | А | MAL | 1LOJ | А | MPD | 1WB8 | А | PMS | 2HQU | С | DUP |
| 1CLS | D | DEC | 1LVW | D | TYD | 1WV0 | А | BN4 | 2HXM | А | 302 |
| 1CPC | А | CYC | 1M2Z | А | BOG | 1WV7 | Т | FUC | 2117 | А | CIT |
| 1CY2 | А | TMP | 1MBQ | А | BEN | 1X29 | В | PMG | 2IPF | А | TRS |
| 1DBN | В | NAG | 1MFI | В | FHC | 1XEY | А | GUA | 2IW6 | А | QQ2 |
| 1DHK | В | NAG | 1MPF | А | C8E | 1XJI | А | D10 | 2IWZ | А | 6NA |
| 1EKX | С | PAL | 1NGP | L | NPA | 1XKD | А | NAP | 2J6E | В | MPD |
| 1EST | А | TSU | 1NIP | В | ADP | 1XR8 | А | PG4 | 2J7L | А | XC2 |
| 1EWY | С | FAD | 1O4H | А | 772 | 1XXS | В | STE | 2J8C | М | GGD |
| 1F42 | А | MNB | 105D | Н | CR9 | 1Y11 | А | 1PE | 2J9C | С | ATP |
| 1FLJ | А | GSH | 106T | А | MES | 1Y2F | А | WAI | 2JH0 | D | 701 |
| 1FQ6 | А | 0QF | 1O9T | В | ATP | 1YRX | В | D9G | 2NY0 | А | HEZ |
| 1G4I | А | MPD | 10AU | J | DNF | 1ZL0 | А | TLA | 20IZ | А | TSR |
| 1G5N | А | SGN | 1PFK | А | ADP | 1ZOM | А | 339 | 20L4 | В | JPN |
| 1G8I | A | P6G | 1Q4J | А | GTX | 1ZRK | А | 367 | 20M9 | А | AJA |
| 1GG6 | С | APF | 1Q6O | В | LG6 | 2A01 | А | AC9 | 20PY | А | CO9 |
| 1GKA | A | D12 | 1Q6Y | А | MPD | 2APX | А | MLA | 2P95 | А | ME5 |
| 1GMR | А | 2GP | 1QIW | A | DPD | 2AY9 | В | 5PV | 2PL7 | В | HTG |
| 1GOY | А | 3GP | 1RE2 | A | NAG | 2AZ5 | В | 307 | 2UUE | В | GVC |
| 1GZR | В | C15 | 1RH7 | С | P6G | 2B0U | В | MPD | 2YXJ | A | N3C |
| 1H1B | A | 151 | 1RHM | В | NA4 | 2B45 | Х | EPE | 35C8 | L | NOX |
| 1H48 | С | CDI | 1RTK | A | GBS | 2BUQ | А | CAQ | 3LJR | А | GGC |
| 1HJ1 | A | PMB | 1RYD | A | GLC | 2C01 | Х | ATP | 4LIP | D | CCP |
| 1HUR | A | GDP | 1RZH | Н | CDL | 2C4L | А | SIA | 4VGC | В | SRD |
| 1HVV | A | TAR | 1S57 | В | EPE | 2C97 | В | MPD | 6RNT | А | 2AM |
| 1HX0 | А | AC1 | 1S9Q | В | CHD | 2CL0 | Х | TRS | 9RSA | В | ADU |
| 115G | A | TS5 | 1SUP | A | PMS | 2CZ5 | В | CIT | | | |
| 119B | D | EPE | 1SVL | С | ADP | 2DCY | A | TAR | | | |
| 1ICR | A | NIO | 1TB6 | 1 | MPD | 2DJH | A | UM3 | | | |

Dataset of PP complexes from the PIBASE database (Table C).

| | PDB ID | Chain | Chain | | | | | | | | | | | | |
|---|---------|--------|--------|------|--------|---|--------------|--------|---|-------|---|---|------|--------|--------|
| | 1YM4 | А | С | 1CA7 | А | С | 1EZR | А | В | 1HAK | А | В | 1M43 | А | В |
| | 1A22 | В | А | 1CBK | А | В | 1EZV | В | А | 1HG4 | А | D | 1M48 | А | В |
| | 1A2K | С | В | 1CD1 | А | В | 1EZV | А | В | 1HKV | А | В | 1M56 | А | В |
| | 1A2V | С | D | 1CD9 | В | А | 1EZV | С | D | 1HP0 | А | в | 1M7N | А | В |
| | 1A5A | A | в | 1CH4 | А | в | 1F7V | D | С | 1HR6 | D | С | 1M7W | А | в |
| | 1454 | R | Δ | 1CKI | Δ | B | 1F1.I | Δ | B | 1HR6 | Δ | F | 1MCZ | Δ | B |
| | 1A7K | C | | | Δ | B | 1520 | Δ | B | | Δ | B | 1ME8 | R | Δ |
| | 1486 | Δ | B | 1CN/ | Δ | C | 1F3T | Ĉ | Б | | R | Δ | 1MG2 | 0 | N |
| | 1/100 | ^ | B | | л л | B | 1E5D | ^ | B | | ^ | Ĉ | | ^ | |
| | 1400 | R | Δ | | Δ | B | 1E60 | Δ | B | 1107 | Δ | B | | Δ | B |
| | 1488 | ^ | R | | л л | B | 1E6B | л л | B | 11102 | ~ | B | 1M71 | ^ | B |
| | | | ы Г | | ^ | C | 1E6D | ^ | C | 1120 | | C | | ~ | C C |
| | | | I D | 100 | A ^ | | 100 | A ^ | | 1144 | | C | | ~ | |
| | | A A | D | 1000 | A A | D | | | D | 1141 | | B | | A ^ | D |
| | IAF3 | A 0 | D T | | A | | | | | 1140 | Ā | D | | A | D |
| | TAGT | 0 | | 1D4A | A | | | A | В | 11/1 | A | В | | A | В |
| | TAGR | A | D | 1065 | A | В | 1FRI 4FOK | A | В | 1185 | Ď | В | 1N40 | A | В |
| | TAIG | P | 0 | 1060 | A | В | 1FSK | J | L | 1185 | A | В | | A | В |
| | 1AIG | L | H | 1080 | A | В | 1FTM | A | C | 1116 | A | В | 1N9E | A | В |
| | 1AIP | A | C | 1DBQ | A | В | 1FUJ | A | D | 1IJL | A | В | 1NBU | В | C |
| | 1AJS | A | В | 1DCL | A | В | | В | A | 11Y8 | в | C | 1NCC | N | Н |
| | 1AR1 | В | A | 1DF8 | A | В | 1FX0 | В | A | 1IYK | A | В | 1NF3 | в | D |
| | 1AUW | C | D | 1DIR | A | В | 1G1A | A | В | 1IZ1 | A | P | 1NFD | A | D |
| | 1AZS | A | В | 1DJU | A | В | 1G2O | A | С | 1J7E | Α | В | 1NFQ | A | В |
| | 1AZS | С | В | 1DN0 | В | Α | 1G4A | В | A | 1J90 | Α | В | 1NKS | A | В |
| | 1AZZ | С | A | 1DO5 | A | С | 1G5Q | A | D | 1JDS | А | С | 1NTO | С | D |
| | 1B3R | В | D | 1DOF | В | С | 1G6O | А | В | 1JEB | А | D | 1NW4 | А | В |
| | 1B6C | D | С | 1DOH | А | В | 1G6Y | А | В | 1JFF | А | В | 1NX9 | В | D |
| | 1B7G | 0 | Q | 1DPJ | А | В | 1G73 | С | А | 1JG8 | А | В | 104S | А | В |
| | 1B7T | Z | А | 1DTY | А | В | 1G85 | А | В | 1JKF | А | В | 105D | L | Т |
| | 1B7T | Y | А | 1DUG | А | В | 1G8Y | С | D | 1JSW | А | D | 1051 | А | В |
| | 1B8G | А | В | 1DVR | А | В | 1G9M | С | G | 1JT0 | А | С | 1050 | А | В |
| | 1B9C | А | В | 1DZB | А | В | 1GEG | В | С | 1JWH | А | D | 1061 | А | В |
| | 1BAI | А | В | 1E3I | А | В | 1GH7 | А | В | 1JWI | В | А | 1063 | А | В |
| | 1BB3 | А | В | 1E3U | А | С | 1GMY | А | В | 1K2O | А | В | 106E | В | А |
| | 1BCC | Е | D | 1E5E | А | В | 1GNX | А | В | 1K3F | А | D | 1O9J | А | В |
| | 1BCC | G | С | 1E7W | А | В | 1GP7 | А | В | 1KBU | А | В | 10AT | В | С |
| | 1BD3 | А | С | 1E8T | А | В | 1GPM | А | В | 1KCZ | А | В | 10D2 | А | В |
| | 1BEB | А | В | 1E9S | Е | F | 1GPW | А | В | 1KFY | В | С | 10DL | А | В |
| | 1BGY | R | 0 | 1ECE | А | В | 1GR7 | А | D | 1KFY | D | В | 10E7 | А | В |
| | 1BI7 | А | В | 1ECS | А | В | 1GRI | А | В | 1KFY | С | В | 10GA | D | Е |
| | 1BJF | А | В | 1EE0 | А | В | 1GTV | А | В | 1KI9 | А | В | 10GA | Е | D |
| | 1BKJ | А | В | 1EI1 | А | В | 1GWN | А | С | 1KOB | А | В | 1OJ6 | А | В |
| | 1BKO | А | В | 1EK6 | А | В | 1GYL | А | В | 1KRU | А | С | 10ME | А | В |
| | 1BMF | С | D | 1EKX | В | С | 1GZ6 | А | В | 1KSG | А | В | 10PL | А | В |
| | 1BML | А | С | 1EM6 | А | В | 1GZM | А | В | 1KXJ | А | В | 10RR | А | В |
| | 1BQH | G | Н | 1EOC | В | А | 1H48 | А | В | 1L8X | А | В | 10RT | С | D |
| | 1BUQ | А | В | 1EP1 | А | в | 1H4G | А | В | 1LBH | А | В | 10RW | В | D |
| | 1BVR | А | В | 1EQU | А | в | 1H5B | А | В | 1LLU | Е | F | 10VL | А | В |
| | 1BYE | С | D | 1EWK | А | в | 1H5Q | А | В | 1LRT | А | В | 10VM | В | D |
| | 1C0T | А | В | 1EYS | М | С | 1H8V | D | Е | 1LW6 | Е | I | 10YJ | С | D |
| ~ | ontinue | 1 | l | - | | - | | | | | | | | - | _ |

| 10ZF | А | В | 1SVX | В | А | 1W6T | А | В | 2BVC | D | Е |
|-------|--------|--------|-------|--------|--------|--------|--------|--------|---------------|--------|---|
| 1P0K | А | В | 1SXG | А | D | 1W6U | А | В | 2C7F | А | D |
| 1P0S | н | Е | 1T2O | в | А | 1W91 | F | G | 2CLT | А | В |
| 1P4B | L | н | 1T8S | А | в | 1WC1 | в | C | 2D1Y | в | С |
| 1P4E | Ā | D | 1T91 | A | D | 1WUE | Ā | B | 2DFT | Ā | C |
| 1P60 | A | B | 1T97 | A | B | 1WYF | A | B | 2DNS | A | B |
| 1P7C | Δ | B | 1TAH | Δ | B | 1\//78 | Δ | B | 2DW6 | Δ | B |
| 103 | Δ | C | 1TC0 | Δ | B | 11/20 | Δ | F | 2D110 | П | B |
| 1040 | ~ | B | | ~ | B | 1727 | - - | ĸ | 2014 2ECH | B | |
| | | B | | | C C | 170 | ٦ ٨ | R | | ^ | |
| | A | D | | | C B | | A D | | | A D | |
| | | D | | A | D | | | A | | | |
| | | D | | A | В | | A | В | | A | В |
| 1PVVX | в | C | 11002 | A | В | 1XG5 | A | В | 2F73 | A | В |
| 1PY4 | A | C | 1 X | A | В | 1XGM | A | В | 2FG6 | C | D |
| 1PYT | В | A | 1TZ3 | A | В | 1XI9 | A | В | 2FM6 | A | В |
| 1Q0C | А | С | 1U0L | A | С | 1XKQ | Α | В | 2FYI | Α | В |
| 1Q3D | А | В | 1U1I | А | В | 1XNX | А | В | 2G2U | А | В |
| 1Q3G | А | В | 1U46 | А | В | 1XSE | А | В | 2GF0 | А | В |
| 1Q57 | А | С | 1U5Q | А | В | 1XTT | В | С | 2GIC | А | В |
| 1Q5Q | С | В | 1UB7 | А | В | 1XXD | В | С | 2GJ7 | Е | В |
| 1Q6T | А | В | 1UFH | В | А | 1XYG | В | С | 216A | А | D |
| 1Q8M | А | В | 1UI5 | А | В | 1Y1M | А | В | 2IFA | А | В |
| 1Q90 | В | D | 1UIM | А | В | 1Y1P | А | В | 2IG2 | Н | L |
| 1QA9 | А | В | 1UIU | А | В | 1YDE | В | Р | 2IGO | С | D |
| 1QPB | А | В | 1UKM | А | В | 1YO6 | А | В | 2IJ2 | В | А |
| 1QZF | С | D | 1UKV | Y | G | 1YP2 | С | D | 2J12 | В | А |
| 1R4P | Е | А | 1UMO | А | В | 1YQ2 | А | В | 2NUU | В | С |
| 1R5K | А | С | 1UPA | в | D | 1YTA | А | В | 2023 | А | В |
| 1R8Q | Е | A | 1URZ | в | С | 1YTZ | С | 1 | 202Y | С | D |
| 1RD5 | А | В | 1USI | А | С | 1YVB | А | 1 | 20KR | А | D |
| 1RD7 | A | B | 1UTR | A | B | 1YXM | A | B | 20NI | С | Ā |
| 1RE5 | B | C | 1000 | Δ | B | 1708 | B | D | 2PD3 | B | C |
| 1R.IN | Δ | B | 1UVQ | Δ | B | 171B | Δ | B | 2PMT | Ċ | D |
| | Δ | D | 1117M | Δ | B | 1754 | Δ | B | 2ΤΔΔ | Δ | B |
| | Δ | B | 1\/2 | Δ | B | 177M | F | Δ | 21/01 3TAT | Ĉ | D |
| | Δ | C | 1\//1 | Δ | B | 1778 | × | Ŵ | | П | C |
| 1920 | ~ | B | 11/08 | ~ | B | | | G | | D | U |
| 1020 | ~ | B | | | ы | | | 5 | | | |
| 1902 | | B | | ^ | B | 12FD | ^ | | | | |
| 1002 | A C | В 7 | | ~ ^ | B | 2420 | A C | | | | |
| 1000 | G | ~ | | A D | | | G ^ | | | | |
| 1510 | | A | | В | A | ZA9K | A | В | | | |
| 15110 | A | | | A | D | ZAYK | В | A | | | |
| 15P8 | A | В | TVIQ | В | C D | | A | C D | | | |
| 1SPG | в | A | 1VIY | A | В | 2AG5 | A | В | | | |
| 1SPI | A | В | 1VKG | A | В | 2ANC | C | E | | | |
| 1SQB | В | I | 1VM6 | В | D | 2B4G | A | В | | | |
| 1SQI | Α | В | 1VPX | С | В | 2BCK | А | D | | | |
| 1SQL | D | E | 1W0M | В | D | 2BD0 | А | В | | | |
| 1SQP | А | Ι | 1W2Z | А | В | 2BEX | С | А | | | |
| 1SU2 | А | В | 1W59 | А | В | 2BTW | А | В | | | |

| Dataset of PL complexes from | the PIBASE database (| (Table D). |
|------------------------------|-----------------------|------------|
|------------------------------|-----------------------|------------|

| PDB ID | Chain | Ligand | | | | | , | | | | |
|---------------|--------|--------|---------------|---------|---------------------|--------------|--------|---------------|--------------|--------|--------|
| 1A07 | A | PTR | 1EE2 | А | CHD | 1LW5 | С | PLG | 1U2R | А | DDE |
| 1A2D | А | PYX | 1EEI | D | GAA | 1M2Z | А | BOG | 1U4S | А | BIH |
| 1A50 | А | FIP | 1EFA | А | NPF | 1MVN | А | PCO | 1U7G | А | BOG |
| 1A54 | А | MDC | 1EHK | А | BNG | 1N0S | А | FLU | 1U9V | А | IHE |
| 1A69 | A | FMB | 1EHK | В | BNG | 1N3I | В | DIH | 1UAY | A | ADN |
| 1A8G | A | 2ZN | 1ELU | Ā | PDA | 1N3Z | Ā | ADN | 1UTR | A | PCB |
| 1A9X | B | CYG | 1FUP | A | ASD | 1NBP | Α | MHC | 1UYS | A | H1I |
| 1AB8 | A | FOK | 1FYN | A | 2AN | 1NI P | C | MN8 | 1VAF | A | ARR |
| 1ACM | A | PAI | 1FYS | M | BGI | 1NX9 | Ā | AIC | 1VFA | A | HBN |
| 1AHI | A | CHO | 1FYS | Н | BGI | 10KF | A | BOG | 1VKG | B | CRI |
| 1AOF | A | GW3 | 1F40 | A | GPI | 10SV | A | CHC | 1W9D | M | SEH |
| | B | EST | 1FL6 | I | | 1PKG | Α | PTR | 1X8 | A | BMH |
| 1AWB | B | | 1FP1 | D | HCC | 1PP9 | Δ | BHG | 1X8\/ | Δ | ESI |
| 1R4W | Δ | BOG | 1FLIP | Δ | PMA | 1PP9 | B | BHG | 1269 | Δ | P34 |
| 1B59 | Δ | | 1691/ | Δ | RO3 | 1PP9 | D | BHG | | Δ | MM1 |
| 1B00 | Δ | CRO | 1665 | Δ | F09 | 1PP9 | F | BHG | 17K3 | Л | BOG |
| 1890 | Δ | PXG | 1665 | Δ | ΕΔΔ | 1001 | G | BHG | 1VP2 | Δ | PMB |
| 189\/ | Δ | RA2 | 1425 | Δ | BOG | 100B | Δ | ΝΔΤ | 177V | Δ | |
| 1BA\/ | Δ | RID | 11/20 | Δ | GPP | 1034 | Ĉ | NGH | 248T | Δ | |
| | F | BOG | 1H4G | Δ | FYP | 100/ | Δ | KHD | 2450 | Δ | CPT |
| | ^ | 270 | 1461 | ^ | אחס | 107P | ~ | | 2410 | R | SC3 |
| | ^ | SB6 | 11101 | ^ | F D N K 2 1 | 10/7 | ~ | | 2AJ0 2BM2 | C | DM2 |
| | ^ | | | ^ | DMR | 10.42 | ~ | | | B | |
| 1DL3 1BM7 | ^ | | | ^ | | 1280 | | | 2014 | ^ | 360 |
| | ^ | | | ^ | | | | | 2003 | | 200 |
| | ^ | | | ^ | BNC | | ~ | | | і Л | 1 |
| | ^ | | 1153 | ^ | PTC | | ~ | | | R | FEE |
| 1BUQ 1BY/ | ^ | | 1101 | ^ | RNI | | R | 550 | | ^ | |
| 1B71 | Δ | | | Δ | CB1 | | Δ | Δ <i>1</i> 5 | 2000 | Δ | |
| 1020 1010 | Δ | RAI | | Δ | SER | 1920 | Δ | 745 FI F | 2037 2EVC | Δ | FC3 |
| 1011 | ^ | | | ^ | | 1920 | ~ | | 200 | ^ | |
| 1000 | ^ | | 1108 | R | 5DA | 1920 | ~ | | | ц | |
| 1C9C | ^ | SEB | 1 112 | ^ | | 1900 | ~ | | | Λ Λ | |
| 1001 | ^ | | | л Ц | ECO | 19570 | ~ | | 2010 | ^ | |
| 1001 | R | | | ۱۱ ۸ | | 1911 | ~ | | 2010 | ^ | |
| 1001 | F | | | ^ | | 1010 | ~ | MBO | | R | |
| 1000 | ۱ ۸ | BOG | 1 IV/N | ^ | 1/3 | | ~ | | 211 1 | ^ | 92D |
| | R | | 1 1 1 1 1 1 1 | Δ | 145 | 1500 | Δ | 860 | 200W | Δ | 014 |
| 1071 | F | | 1K3T | Ċ | BR7 | 1901 | Δ | BTS | 2001 | Δ | C 5 3 |
| 1021 1D3H | | Δ26 | | Δ | | 1907 1978 | Δ | PIG | 20Q1 | R | |
| 1D311 1D4E | Δ | | | Ċ | BRS | 102IX | R | IMR | 2011 | Δ | |
| | ^ | | | D D | BDC | 1T2T | ^ | | 2011 | ^ | |
| | Δ | | | Δ | | 1T6H | Δ | | 3001 | Δ | |
| | ^ | GNT | 11/10 | ^ | MC | 1783 | ~ | | | ^ | СМЦ |
| 1070 | ^ | 97V | | ^ | DOT | 1105 1T0B | ~ | 102 | 4703 6000 | ī | CIVILI |
| 1E2N | Δ | | | R | 1/7 | | ~ | | 1100 | ∟ ∧ | |
| | R | E36 | | ت ۸ | י <i>ייו</i> עדס | | ~ | | 1103 | л | EAA |
| | F | MQ1 | | л С | FMY | | ~ | FMR | | | |
| 101 | | | | | | | ~ | | | | |
| | D | | | | | | A D | 330 2 V VI | | | |
| IECZ | D | DUG | ILKU | А | BPO | | В | ∠AN | 1 | | |

| PDB ID | Chain | Ligand | | | | | | |
|--------|-------|--------|------|---|-----|------|---|-----|
| 1A6A | А | NAG | 2LZG | А | 13Q | 3M1I | А | GTP |
| 1A7X | В | FKA | 20U7 | А | ANP | 3MDY | А | LDN |
| 1AB8 | А | FOK | 2P1M | В | IHP | 3NC0 | А | IPH |
| 1AGR | А | GDP | 2P26 | А | NAG | 3047 | А | GDP |
| 1AZS | С | GSP | 2TNF | А | TRS | 3OSK | А | NAG |
| 1CS4 | А | MES | 2VJE | D | FLC | 3QAZ | В | NAG |
| 1D8D | В | FII | 2VOH | А | CIT | 3QBR | Х | NHE |
| 1H2K | А | OGA | 2W5Y | А | SAH | 3QD6 | А | NAG |
| 118L | А | NAG | 2WBE | С | ANP | 3QTK | А | TFA |
| 1IA0 | K | ACP | 2WKP | А | FMN | 3SL9 | А | IMD |
| 1ICF | I | NAG | 2XRP | В | GTP | 3TX7 | В | P6L |
| 1L2I | А | ETC | 2YDS | А | NAG | 3U88 | А | GGB |
| 109C | А | FLC | 2YEM | А | WSH | 3UP0 | А | D7S |
| 1P93 | А | DEX | 2ZKW | А | CU1 | 3V4P | С | TRS |
| 1PZN | В | IMD | 3A9E | В | REA | 3VNG | А | FUU |
| 1QAB | Е | RTL | 3BEJ | А | MUF | 4A9E | А | 3PF |
| 1R6N | А | 434 | 3CLX | D | X22 | 4AY6 | А | 12V |
| 1S1J | А | IQZ | 3DZU | D | PLB | 4DBB | А | ACY |
| 1TCO | В | MYR | 3DZU | А | REA | 4DEW | А | LU2 |
| 1TW6 | В | BTB | 3F7Q | А | 1PE | 4DSN | А | GCP |
| 1U27 | А | 4IP | 3FCS | А | IMD | 4E2T | А | EPE |
| 1V7P | В | NAG | 3GBG | А | PAM | 4G1E | А | NAG |
| 1XLS | E | TCD | 3HOF | В | DHC | 4G1M | А | NAG |
| 1YJD | С | NAG | 3HQR | А | OGA | 4GMX | С | K85 |
| 1YSG | А | 4FC | 3HVL | А | SRL | 4GS6 | А | 1FM |
| 2BRQ | А | GSH | 3ICI | А | MES | 4GZ9 | А | NAG |
| 2ERJ | В | NAG | 3IPQ | А | 965 | 4H71 | А | PXE |
| 2GXA | А | ADP | 3IT8 | D | NAG | 4JWL | А | HRC |
| 2H61 | Е | PG4 | 3K6S | А | MAN | 4LOO | А | SB4 |
| 2HBH | А | XE4 | 3L0L | А | HC3 | | | |

Dataset of PL complexes from the TIMBAL database (Table E).

| No | PDB ID | Resolution (Å) | Name of homodimer | Scientific source | Chain one | Length | Chain two | Length |
|----|-----------|-------------------|---|--|--------------|--------|--------------|--------|
| 1 | 1GIF | 1.90 | Human Glycosylation-Inhibiting Factor | Homo sapiens | В | 115 | А | 15 |
| 2 | 1A09 | 2.00 | Peptide Ligands Of PP60(C-SRC) SH2 Domains | Homo sapiens | А | 107 | В | 107 |
| 3 | 1AB8 | 2.20 | Adenylyl Cyclase | Homo sapiens | В | 220 | А | 220 |
| 4 | 1ASL | 2.60 | Aspartate Aminotransferase | Escherichia coli | В | 396 | А | 396 |
| 5 | 1AV1 | 4.00 | Apolipoprotein A-I | Homo sapiens | В | 201 | А | 201 |
| 6 | 1B06 | 2.20 | Protein (Superoxide Dismutase) | Sulfolobus acidocaldarius DSM 639 | А | 210 | В | 210 |
| 7 | 1B99 | 2.70 | Protein (Nucleoside Diphosphate Kinase) | Dictyostelium discoideum | F | 155 | С | 155 |
| 8 | 1BJQ | 2.65 | Lectin | Vigna unguiculata subsp. cylindrica (sow-pea) | С | 253 | А | 253 |
| 9 | 1BRR | 2.90 | Protein (Bacteriorhodopsin) | Halobacterium salinarum | А | 247 | С | 247 |
| 10 | 1BVI | 1.90 | Protein (Ribonuclease T1) | Aspergillus oryzae | А | 104 | С | 104 |
| 11 | 1C2O | 4.20 | Acetylcholinesterase | Electrophorus electricus | D | 539 | А | 539 |
| 12 | 1CKG | 2.20 | Protein (Lysozyme) | Homo sapiens | В | 130 | А | 130 |
| 13 | 1CYY | 2.15 | DNA Topoisomerase I | Escherichia coli K-12 | А | 264 | В | 264 |
| 14 | 1D7F | 1.90 | Cyclodextrin Glucanotransferase | Bacillus sp. 1011 | В | 686 | А | 686 |
| 15 | 1DLE | 2.10 | Complement Factor B | Homo sapiens | А | 298 | В | 298 |
| 16 | 1DM5 | 1.93 | Annexin XII E105k Mutant Homohexamer | Hydra vulgaris | В | 315 | D | 315 |
| 17 | 1E8U | 2.00 | Hemagglutinin-Neuraminidase | Newcastle disease virus | В | 454 | А | 454 |
| 18 | 1EZL | 2.00 | Azurin | Pseudomonas aeruginosa | С | 128 | D | 128 |
| 19 | 1F46 | 1.50 | Cell Division Protein Zipa | Escherichia coli | В | 140 | А | 140 |
| 20 | 1F88 | 2.80 | Rhodopsin | Bos taurus | В | 348 | А | 348 |
| 21 | 1FX9 | 2.00 | Phospholipase A2, Major Isoenzyme | Sus scrofa | В | 124 | А | 124 |
| 22 | 1G0A | 2.04 | Hemoglobin Beta Chain | Bos taurus | D | 145 | В | 145 |
| 23 | 1G21 | 3.00 | Nitrogenase Iron Protein | Azotobacter vinelandii | Н | 289 | G | 289 |

Table 3.2 Dataset of PP complexes from the ABC datasetList of PP complexes from the ABC dataset comprised of 94 homodimer complexes (Table A).

| 24 | 1GQ3 | 2.01 | Aspartate Carbamoyltransferase | Escherichia coli K-12 | С | 310 | В | 310 |
|----|------|------|--|----------------------------|---|-----|---|-----|
| 25 | 1H6A | 2.50 | Precursor Form Of Glucose-Fructose Oxidoreductase | Zymomonas mobilis | В | 433 | А | 433 |
| 26 | 1HFY | 2.30 | Alpha-Lactalbumin | Capra hircus | А | 123 | В | 123 |
| 27 | 1HQK | 1.60 | 6,7-Dimethyl-8-Ribityllumazine Synthase | Aquifex aeolicus | Е | 154 | D | 154 |
| 28 | 1HWU | 2.10 | PII Protein | Herbaspirillum seropedicae | В | 112 | А | 112 |
| 29 | 118F | 1.75 | Putative Snrnp SM-Like Protein | Pyrobaculum aerophilum | С | 81 | В | 81 |
| 30 | 1IFV | 2.25 | Protein LLR18B | Lupinus luteus | А | 155 | В | 155 |
| 31 | 1IIN | 2.10 | Glucose-1-Phosphate Thymidylyltransferase | Salmonella enterica | А | 292 | В | 292 |
| 32 | 1KQD | 1.90 | Oxygen-Insensitive NAD(P)H Nitroreductase | Enterobacter cloacae | А | 217 | В | 217 |
| 33 | 1KWS | 2.10 | Beta-1,3-Glucuronyltransferase 3 | Homo sapiens | В | 261 | А | 261 |
| 34 | 1MCI | 2.70 | Immunoglobulin Lambda Dimer MCG (Light Chain) | Homo sapiens | В | 216 | А | 216 |
| 35 | 1MI3 | 1.80 | Xylose Reductase | Candida tenuis | В | 319 | А | 319 |
| 36 | 1MSA | 2.29 | Agglutinin | Galanthus nivalis | В | 109 | С | 109 |
| 37 | 1N9S | 3.50 | Small Nuclear Ribonucleoprotein F | Saccharomyces cerevisiae | С | 93 | D | 93 |
| 38 | 1NYS | 3.05 | Activin Receptor | Rattus norvegicus | А | 105 | С | 105 |
| 39 | 1081 | 1.50 | Tryparedoxin II | Crithidia fasciculata | В | 152 | А | 152 |
| 40 | 10BQ | 1.85 | Crustacyanin C1 Subunit | Homarus gammarus | В | 181 | А | 181 |
| 41 | 10KT | 1.90 | Glutathione S-Transferase | Plasmodium falciparum | В | 211 | А | 211 |
| 42 | 10ME | 2.30 | Beta-Lactamase | Staphylococcus aureus | А | 258 | В | 258 |
| 43 | 10PF | 3.20 | Matrix Porin Outer Membrane Protein F | Escherichia coli | В | 340 | А | 340 |
| 44 | 1P7L | 2.50 | S-Adenosylmethionine Synthetase | EScherichia coli | А | 383 | В | 383 |
| 45 | 1PA0 | 2.20 | Myotoxic Phospholipase A2-Like | Bothrops pauloensis | В | 121 | А | 121 |
| 46 | 1PMO | 2.30 | Glutamate Decarboxylase Beta | Escherichia coli A | D | 466 | С | 466 |
| | | | | | | | | |

| 47 | 1PYG | 2.87 | Glycogen Phosphorylase B | Oryctolagus cuniculus | В | 842 | А | 842 |
|----|------|------|---|---|---|-----|---|-----|
| 48 | 1QPX | 2.40 | Papd Chaperone | Escherichia coli | А | 218 | В | 218 |
| 49 | 1R5T | 2.00 | Cytidine Deaminase | Saccharomyces cerevisiae | А | 142 | В | 142 |
| 50 | 1RFX | 2.00 | Resistin | Mus musculus | В | 94 | С | 94 |
| 51 | 1RH7 | 3.11 | Resistin-Like Beta | Mus musculus | С | 81 | В | 81 |
| 52 | 1RQ7 | 2.60 | Cell Division Protein FTSZ | Mycobacterium tuberculosis | В | 382 | А | 382 |
| 53 | 1S7Y | 1.75 | Glutamate Receptor, Ionotropic Kainate 2 Precursor | Rattus norvegicus | В | 259 | А | 259 |
| 54 | 1SND | 1.84 | Staphylococcal Nuclease Dimer | Staphylococcus aureus | А | 143 | В | 143 |
| 55 | 1TH7 | 1.68 | Small Nuclear Riboprotein Protein | Sulfolobus solfataricus | L | 81 | K | 81 |
| 56 | 1U3R | 2.21 | Estrogen Receptor Beta | Homo sapiens | А | 241 | В | 241 |
| 57 | 1U3W | 1.45 | Alcohol Dehydrogenase Gamma Chain | Homo sapiens | В | 374 | А | 374 |
| 58 | 1UVC | 2.00 | Nonspecific Lipid Transfer Protein | Oryza sativa | А | 91 | В | 91 |
| 59 | 1VGQ | 2.13 | Formyl-Coenzyme A Transferase | Oxalobacter formigenes | В | 427 | А | 427 |
| 60 | 1VLZ | 2.05 | Chey | Escherichia coli | В | 128 | А | 128 |
| 61 | 1W0I | 2.10 | 3-Oxoacyl Carrier Protein Synthase | Arabidopsis thaliana | В | 431 | А | 431 |
| 62 | 1W29 | 2.30 | 6,7-Dimethyl-8-Ribityllumazine Synthase | Mycobacterium tuberculosis | С | 160 | D | 160 |
| 63 | 1X1Z | 1.45 | Orotidine 5'-Phosphate Decarboxylase | Methanothermobacter thermautotrophicus | А | 252 | В | 252 |
| 64 | 1XBY | 1.58 | 3-Keto-L-Gulonate 6-Phosphate Decarboxylase | Escherichia coli | В | 216 | А | 216 |
| 65 | 1XMZ | 1.38 | GFP-Like Chromoprotein FP595 | Anemonia sulcata | В | 243 | А | 243 |
| 66 | 1YGP | 2.80 | Yeast Glycogen Phosphorylase | Saccharomyces cerevisiae | А | 879 | В | 879 |
| 67 | 1ZRS | 1.50 | Hypothetical Protein | Pseudomonas aeruginosa | В | 317 | А | 317 |
| 68 | 1ZUX | 1.85 | Green To Red Photoconvertible GPF-Like Protein EOSFP | Lobophyllia hemprichii | D | 226 | В | 226 |
| 69 | 2AAX | 1.75 | Mineralocorticoid Receptor | Homo sapiens | А | 275 | В | 275 |
| 70 | 2AMT | 2.30 | 2-C-Methyl-D-Erythritol 2,4- Cyclodiphosphate Synthase | EScherichia coli | С | 159 | В | 159 |

| 71 | 2AQ3 | 2.30 | T-Cell Receptor Beta Chain V | Mus musculus | G | 112 | А | 112 |
|----|------|------|---|---|---|-----|---|-----|
| 72 | 2AY5 | 2.40 | Aromatic Amino Acid Aminotransferase | Paracoccus denitrificans | А | 394 | В | 394 |
| 73 | 2AYQ | 3.00 | 3-Isopropylmalate Dehydrogenase | Bacillus coagulans | В | 366 | А | 366 |
| 74 | 2C3N | 1.50 | Glutathione S-Transferase Theta 1 | Homo sapiens | А | 247 | В | 247 |
| 75 | 2CHP | 2.00 | Metalloregulation DNA-Binding Stress Protein | Bacillus subtilis subsp. subtilis str. 168 | С | 153 | D | 153 |
| 76 | 2COG | 2.10 | branched chain aminotransferase 1, cytosolic | Homo sapiens | А | 386 | В | 386 |
| 77 | 2CUY | 2.10 | Malonyl CoA-[acyl carrier protein] transacylase | Thermus thermophilus HB8 | В | 305 | А | 305 |
| 78 | 2D4V | 1.90 | Isocitrate Dehydrogenase | isocitrate dehydrogenase | В | 429 | А | 429 |
| 79 | 2DVG | 2.78 | Galactose-Binding Lectin | Arachis hypogaea | С | 236 | В | 236 |
| 80 | 2EV4 | 2.28 | Hypothetical Protein Rv1264/MT1302 | Mycobacterium tuberculosis | В | 222 | А | 222 |
| 81 | 2FR6 | 2.07 | Cytidine Deaminase | Mus musculus | С | 146 | D | 146 |
| 82 | 2GPV | 2.85 | Estrogen-Related Receptor Gamma | Homo sapiens | D | 230 | В | 230 |
| 83 | 2GQD | 2.30 | 3-Oxoacyl-[Acyl-Carrier-Protein] Synthase 2 | Staphylococcus aureus | В | 437 | А | 437 |
| 84 | 2GVM | 2.30 | Hydrophobin-1 | Trichoderma reesei | В | 75 | А | 75 |
| 85 | 2H1L | 3.16 | Large T Antigen | Simian virus 40 | Е | 370 | F | 370 |
| 86 | 2IPJ | 1.80 | Aldo-Keto Reductase Family 1 Member C2 | Homo sapiens | В | 321 | А | 321 |
| 87 | 2IYG | 2.30 | Appa, Antirepressor Of PPSR, Sensor Of Blue Light | Rhodobacter sphaeroides | В | 124 | А | 124 |
| 88 | 2J3K | 2.80 | NADP-Dependent Oxidoreductase P1 | Arabidopsis thaliana | А | 345 | В | 345 |
| 89 | 2J96 | 2.25 | Phycoerythrocyanin Alpha Chain | Mastigocladus laminosus | В | 162 | А | 162 |
| 90 | 20L1 | 1.80 | Deoxyuridine 5'-Triphosphate Nucleotidohydrolase | Vaccinia virus | А | 147 | С | 147 |
| 91 | 2TNF | 1.40 | Protein (Tumor Necrosis Factor Alpha) | Mus musculus | В | 156 | А | 156 |
| 92 | 4AKE | 2.20 | Adenylate Kinase | Escherichia coli | А | 214 | В | 214 |
| 93 | 6PFK | 2.60 | Phosphofructokinase | Geobacillus stearothermophilus | С | 319 | D | 319 |
| 94 | 7TIM | 1.90 | Triosephosphate Isomerase | Saccharomyces cerevisiae | В | 247 | А | 247 |

| No | PDB ID | Resolution (Å) | Chain one | Name of chain one | Length | Chain two | Name of chain two | Length |
|----|-----------|-------------------|--------------|--|--------|--------------|---|--------|
| 1 | 15C8 | 2.50 | L | IGG 5C8 FAB (Light Chain) | 213 | Н | IGG 5C8 FAB (Heavy Chain) | 217 |
| 2 | 1A6U | 2.10 | Н | B1-8 Fv (Heavy Chain) | 120 | L | B1-8 FV (Light Chain) | 108 |
| 3 | 1AHW | 3.00 | F | Tissue Factor | 219 | А | Immunoglobulin FAB 5G9 (Light Chain) | 214 |
| 4 | 1AVA | 1.90 | А | Barley Alpha-Amylase 2 | 403 | С | Barley Alpha-Amylase/Subtilisin Inhibitor | 181 |
| 5 | 1AY7 | 1.70 | А | Guanyl-Specific Ribonuclease SA | 96 | В | Barstar | 89 |
| 6 | 1AZZ | 2.30 | А | Collagenase | 226 | С | Ecotin | 142 |
| 7 | 1BMF | 2.85 | F | Bovine Mitochondrial F1-Atpase | 482 | В | Bovine Mitochondrial F1-Atpase | 510 |
| 8 | 1BVN | 2.50 | Р | Alpha-amylase | 496 | Т | Tendamistat | 71 |
| 9 | 1BZX | 2.10 | E | Protein (Trypsin) | 222 | I | Protein (Bovine Pancreatic Trypsin Inhibitor) | 58 |
| 10 | 1DFJ | 2.50 | E | Ribonuclease A | 124 | I | Ribonuclease Inhibitor | 456 |
| 11 | 1DPJ | 1.80 | А | Proteinase A | 329 | В | Proteinase Inhibitor IA3 Peptide | 33 |
| 12 | 1E0F | 3.10 | E | Thrombin | 259 | J | Haemadin | 57 |
| 13 | 1EYS | 2.20 | Н | Photosynthetic Reaction Center | 259 | Μ | Photosynthetic Reaction Center | 324 |
| 14 | 1F45 | 2.80 | А | Interleukin-12 Beta Chain | 306 | В | Interleukin-12 Alpha Chain | 197 |
| 15 | 1FAK | 2.10 | Н | Protein (Blood Coagulation Factor Viia) | 254 | I | Protein (5l15) | 55 |
| 16 | 1FI8 | 2.20 | А | Natural Killer Cell Protease 1 | 228 | С | Ecotin | 84 |
| 17 | 1G3N | 2.90 | А | Cyclin-Dependent Kinase 6 | 326 | В | Cyclin-Dependent Kinase 6 Inhibitor | 168 |
| 18 | 1G6V | 3.50 | А | Carbonic Anhydrase | 260 | K | Antibody Heavy Chain | 126 |
| 19 | 1GXS | 2.30 | D | P-(S)-Hydroxymandelonitrile Lyase Chain B | 158 | С | P-(S)-Hydroxymandelonitrile Lyase Chain A | 270 |
| 20 | 1HE8 | 3.00 | В | Transforming Protein P21/H-Ras-1 | 166 | А | Phosphatidylinositol 3-Kinase Catalytic Subunit, Gamma Isoform | 965 |
| 21 | 1HJA | 2.30 | В | Alpha-Chymotrypsin | 131 | I | Ovomucoid Inhibitor | 51 |
| 22 | 1I3O | 2.70 | В | Caspase 3 | 110 | Е | Baculoviral Iap Repeat-Containing Protein 4 | 121 |
| 23 | 1ICF | 2.00 | А | Protein (Cathepsin L: Heavy Chain) | 175 | I | Protein (Invariant Chain) | 65 |

List of PP complexes from the ABC dataset comprised of heterodimer complexes (Table B).

| 24 | 1JSU | 2.30 | В | Cyclin A | 260 | С | P27 | 84 |
|----|------|------|---|---|-----|---|--|-----|
| 25 | 1JTH | 2.00 | D | Syntaxin 1a | 77 | А | SNAP25 | 82 |
| 26 | 1JZD | 2.30 | В | Thiol:Disulfide Interchange Protein DSBC | 220 | С | Thiol:Disulfide Interchange Protein dsbd | 132 |
| 27 | 1KQM | 3.00 | А | Myosin Heavy Chain | 835 | С | Myosin Essential Light Chain | 156 |
| 28 | 1MCV | 1.80 | А | Elastase 1 | 240 | I | HEI-TOE I | 28 |
| 29 | 1N8O | 2.00 | С | Chymotrypsin A, C Chain | 97 | Е | Ecotin | 142 |
| 30 | 1NHG | 2.43 | В | Enoyl-Acyl Carrier Reductase | 229 | D | Enoyl-Acyl Carrier Reductase | 60 |
| 31 | 1NMA | 3.00 | Ν | N9 Neuraminidase | 388 | Н | Fab Nc10 | 122 |
| 32 | 1NW9 | 2.40 | В | Catalytic Domain Of Caspase-9 | 238 | А | A Inhibitor Of Apoptosis Protein 3 | 91 |
| 33 | 106S | 1.80 | А | Internalin A | 461 | В | E-Cadherin | 105 |
| 34 | 1009 | Х | А | Stromelysin-1 | 168 | В | Metalloproteinase Inhibitor 1 | 128 |
| 35 | 1P0S | 2.80 | Н | Coagulation Factor X Precursor | 254 | Е | Ecotin Precursor | 142 |
| 36 | 1P4I | 2.80 | L | Antibody Variable Light Chain | 135 | Н | Antibody Variable Light Chain | 124 |
| 37 | 1PPF | 1.80 | E | Leukocyte Elastase | 218 | I | Ovomucoid Inhibitor | 56 |
| 38 | 1PST | 3.00 | М | Photosynthetic Reaction Center | 296 | Н | Photosynthetic Reaction Center | 237 |
| 39 | 1QGE | 1.70 | D | Protein (Triacylglycerol Hydrolase) | 222 | Е | Protein (Triacylglycerol Hydrolase) | 97 |
| 40 | 1QRN | 2.80 | D | T-Cell Receptor, Alpha Chain | 200 | Е | T-Cell Receptor, Beta Chain | 243 |
| 41 | 1R8Q | 1.86 | А | ADP-Ribosylation Factor 1 | 181 | Е | Arno | 203 |
| 42 | 1S4Y | 2.30 | В | Inhibin Beta A Chain | 116 | А | Activin Receptor Type IIB Precursor | 98 |
| 43 | 1SC1 | 2.60 | А | Interleukin-1 Beta Convertase | 178 | В | Interleukin-1 Beta Convertase | 88 |
| 44 | 1SGF | 3.15 | G | Nerve Growth Factor | 237 | В | Nerve Growth Factor | 118 |
| 45 | 1SGR | 1.80 | E | Streptomyces Griseus Proteinase B | 185 | I | Turkey Ovomucoid Inhibitor | 51 |
| 46 | 1STF | 2.40 | E | Papain | 212 | I | Stefin B | 98 |

| 47 | 1SUV | 1.75 | D | Serotransferrin, N-Lobe | 329 | F | Serotransferrin, C-Lobe | 345 |
|-------|------|------|---|---|-----|---|--|-----|
| 48 | 1T6G | 1.80 | С | Endo-1,4-Beta-Xylanase I | 184 | А | Xylanase Inhibitor | 381 |
| 49 | 1UGH | 1.90 | E | Protein (Uracil-Dna Glycosylase) | 223 | 1 | Protein (Uracil-Dna Glycosylase Inhibitor) | 82 |
| 50 | 1X1U | 2.30 | А | Ribonuclease | 110 | D | Barstar | 89 |
| 51 | 1XX9 | 2.20 | В | Coagulation Factor XI | 238 | D | Ecotin | 142 |
| 52 | 1Y48 | 1.84 | Е | Subtilisin BPN' | 281 | I | Chymotrypsin Inhibitor 2 | 64 |
| 53 | 1YI5 | 4.20 | С | Acetylcholine-Binding Protein | 210 | Н | Long Neurotoxin 1 | 71 |
| 54 | 1YRQ | 2.10 | А | Periplasmic [NiFe] Hydrogenase Small Subunit | 264 | Н | Periplasmic [Nife] Hydrogenase Large Subunit | 549 |
| 55 | 1ZJD | 2.60 | А | Catalytic Domain of Coagulation Factor XI | 237 | В | Kunitz Protease Inhibitory Domain of Protease Nexin II | 57 |
| 56 | 2AGY | 1.10 | В | Aromatic Amine Dehydrogenase | 361 | D | Aromatic Amine Dehydrogenase | 135 |
| 57 | 2B42 | 2.50 | В | Endo-1,4-beta-xylanase A | 185 | A | Xylanase Inhibitor-I | 381 |
| 58 | 2BE6 | 2.50 | А | Calmodulin 2 | 150 | D | Voltage-Dependent L-Type Calcium Channel Alpha-1C Subunit | 37 |
| 59 | 2BEX | 1.99 | С | Nonsecretory Ribonuclease | 135 | А | Ribonuclease Inhibitor | 460 |
| 60 | 2BL0 | 1.75 | В | Myosin Regulatory Light Chain | 145 | А | Major Plasmodial Myosin Heavy Chain | 63 |
| 61 | 2BTO | 2.50 | Т | Thioredoxin 1 | 108 | А | Tubulin Btuba | 473 |
| 62 | 2D26 | 3.20 | А | Alpha-1-Antitrypsin | 358 | В | Alpha-1-Antitrypsin | 36 |
| 63 | 2DSQ | 2.80 | I | Insulin-Like Growth Factor IB | 70 | G | Insulin-Like Growth Factor-Binding Protein 1 | 94 |
| 64 | 2FHZ | 1.15 | А | Colicin-E5 Immunity Protein | 109 | В | Colicin-E5 | 108 |
| 65 | 2G2U | 1.60 | A | Beta-lactamase SHV-1 | 265 | В | Beta-lactamase inhibitory protein | 165 |
| 66 | 2GD4 | 3.30 | I | Antithrombin-III | 443 | Н | Coagulation Factor, Stuart Factor, Stuart- Prower Factor, Contains: Factor X Light Chain; Factor X Heavy Chain; Activated Factor Xa Heavy Chain | 241 |
| 67 | 2GMR | 2.50 | L | Photosynthetic Reaction Center Protein L Chain | 281 | Μ | Photosynthetic Reaction Center Protein M Chain | 307 |
| 68 | 2HI7 | 3.70 | А | Thiol:Disulfide Interchange Protein dsbA | 189 | Н | Disulfide Bond Formation Protein B | 176 |
| 69 | 212R | 3.35 | Н | Kv Channel-Interacting Protein 1 | 180 | D | Potassium Voltage-Gated Channel Subfamily D Member 3 | 144 |
| 70 | 2IDO | 2.10 | А | DNA Polymerase III Epsilon Subunit | 186 | С | Hot Protein | 83 |
| Conti | nued | | | • | | | | |

| 2IJO | 2.30 | А | Polyprotein | 58 | I | Pancreatic Trypsin Inhibitor | 58 |
|------|--|--|---|--|---|--|--|
| 2IWG | 2.35 | D | IG Gamma-1 Chain C | 207 | Е | 52 Kda Ro Protein | 181 |
| 2NX5 | 2.70 | А | HLA-B35 | 276 | D | ELS4 TCR Alpha Chain | 188 |
| 2NY7 | 2.30 | G | Envelope Glycoprotein Gp120 Serine/Threonine-Protein | 317 | Н | Antibody B12, Heavy Chain | 230 |
| 208A | 2.61 | А | Phosphatase Pp1-Gamma Catalytic Subunit | 329 | I | Protein phosphatase inhibitor 2 | 206 |
| 20CC | 2.30 | А | Cytochrome C Oxidase | 514 | В | Cytochrome C Oxidase | 227 |
| 2P1L | 2.50 | А | Apoptosis Regulator Bcl-X | 153 | В | Beclin 1 | 31 |
| 2PRG | 2.30 | А | Peroxisome Proliferator Activated Receptor Gamma | 271 | С | Nuclear Receptor Coactivator SRC-1 | 88 |
| 2PVO | 3.40 | D | Ferredoxin-1 | 96 | А | Ferredoxin-thioredoxin reductase, catalytic chain | 110 |
| 3PCD | 2.10 | А | Protocatechuate 3,4-Dioxygenase | 200 | Μ | Protocatechuate 3,4-Dioxygenase | 238 |
| | 2IJO 2IWG 2NX5 2NY7 2O8A 2OCC 2P1L 2PRG 2PVO 3PCD | 2IJO2.302IWG2.352NX52.702NY72.302O8A2.612OCC2.302P1L2.502PRG2.302PVO3.403PCD2.10 | 2IJO 2.30 A 2IWG 2.35 D 2NX5 2.70 A 2NY7 2.30 G 2O8A 2.61 A 2OCC 2.30 A 2P1L 2.50 A 2PRG 2.30 A | 2IJO2.30APolyprotein2IWG2.35DIG Gamma-1 Chain C2NX52.70AHLA-B352NY72.30GEnvelope Glycoprotein Gp120 Serine/Threonine-Protein2O8A2.61APhosphatase Pp1-Gamma Catalytic Subunit2OCC2.30ACytochrome C Oxidase2P1L2.50AApoptosis Regulator Bcl-X2PRG2.30APeroxisome Proliferator Activated Receptor Gamma2PVO3.40DFerredoxin-13PCD2.10AProtocatechuate 3,4-Dioxygenase | 2IJO2.30APolyprotein582IWG2.35DIG Gamma-1 Chain C2072NX52.70AHLA-B352762NY72.30GEnvelope Glycoprotein Gp1203172O8A2.61APhosphatase Pp1-Gamma Catalytic Subunit3292OCC2.30ACytochrome C Oxidase5142PRG2.30APeroxisome Proliferator Activated Receptor Gamma2712PVO3.40DFerredoxin-1963PCD2.10AProtocatechuate 3,4-Dioxygenase200 | 2IJO2.30APolyprotein58I2IWG2.35DIG Gamma-1 Chain C207E2NX52.70AHLA-B35276D2NY72.30GEnvelope Glycoprotein Gp120317H2O8A2.61APhosphatase Pp1-Gamma Catalytic Subunit329I2OCC2.30ACytochrome C Oxidase514B2P1L2.50AApoptosis Regulator Bcl-X153B2PRG2.30APeroxisome Proliferator Activated Receptor Gamma271C2PVO3.40DFerredoxin-196A3PCD2.10AProtocatechuate 3,4-Dioxygenase200M | 2IJO2.30APolyprotein58IPancreatic Trypsin Inhibitor2IWG2.35DIG Gamma-1 Chain C207E52 Kda Ro Protein2NX52.70AHLA-B35276DELS4 TCR Alpha Chain2NY72.30GEnvelope Glycoprotein Gp120317HAntibody B12, Heavy Chain2O8A2.61APhosphatase Pp1-Gamma Catalytic Subunit329IProtein phosphatase inhibitor 22OCC2.30ACytochrome C Oxidase514BCytochrome C Oxidase2P1L2.50AApoptosis Regulator BcI-X153BBeclin 12PRG2.30APeroxisome Proliferator Activated Receptor Gamma271CNuclear Receptor Coactivator SRC-12PVO3.40DFerredoxin-196AFerredoxin-thioredoxin reductase, catalytic chainFerredoxin-thioredoxin sequences3PCD2.10AProtocatechuate 3,4-Dioxygenase200MProtocatechuate 3,4-Dioxygenase |

| No | PDB ID | Resolution (Å) | Name of homodimer | Scientific source | Chain one | Length | Chain two | Length |
|----|-----------|-------------------|--|-------------------------------|--------------|--------|--------------|--------|
| 1 | 1YM4 | 2.25 | Beta-Secretase 1 | Homo sapiens | А | 408 | С | 408 |
| 2 | 1A2V | 2.40 | Methylamine Oxidase | Ogataea angusta | С | 655 | D | 655 |
| 3 | 1A7K | 2.80 | Glyceraldehyde-3-Phosphate Dehydrogenase | Leishmania mexicana | С | 360 | D | 360 |
| 4 | 1A8G | 2.50 | HIV-1 Protease | Human immunodefiency virus | А | 99 | В | 99 |
| 5 | 1A8T | 2.55 | Metallo-Beta-Lactamase | Bacteroides fragilis | А | 232 | В | 232 |
| 6 | 1A99 | 2.20 | Putrescine-Binding Protein | Escherichia coli | В | 344 | А | 344 |
| 7 | 1AB8 | 2.20 | Adenylyl Cyclase | Rattus norvegia | А | 220 | В | 220 |
| 8 | 1AE1 | 2.40 | Tropinone Reductase-1 | Datura stramonium | А | 273 | В | 273 |
| 9 | 1AFS | 2.50 | 3-Alpha-Hydroxysteroid Dehydrogenase | Rattus norvegia | А | 323 | В | 323 |
| 10 | 1AG1 | 2.36 | Triosephosphate Isomerase | Trypanosoma brucei | 0 | 250 | Т | 250 |
| 11 | 1AGR | 2.80 | Guanine Nucleotide-Binding Protein G(I) | Rattus norvegia | А | 353 | D | 353 |
| 12 | 1AJS | 1.60 | Aspartate Aminotransferase | Sus scrofa | А | 412 | В | 412 |
| 13 | 1AUW | 2.50 | Delta 2 Crytallin | Anas platyrhynchos | С | 468 | D | 468 |
| 14 | 1B3R | 2.80 | Protein (S-Adenosylhomocysteine Hydrolase) | Rattus norvegia | В | 431 | D | 431 |
| 15 | 1B7G | 2.05 | Protein (Glyceraldehyde-3-phosphate dehydrogenase) | Sulfolobus solfatarius | 0 | 340 | Q | 340 |
| 16 | 1B8G | 2.37 | Protein (1-Aminocyclopropane-1 carboxylate synthase) | Malus domestica | А | 429 | В | 429 |
| 17 | 1B9C | 2.40 | Protein (Green Fluorescent Protein) | Auquorea victoria | А | 238 | В | 238 |
| 18 | 1BAI | 2.40 | Protease | Rous sarcoma virus | А | 124 | В | 124 |
| 19 | 1BB3 | 1.80 | Lysozyme | Homo sapiens | А | 130 | В | 130 |
| 20 | 1BD3 | 1.93 | Uracil Phosphoribosyltransferase | Toxoplasma gondii | А | 243 | С | 243 |
| 21 | 1BEB | 1.80 | Beta-Lactoglobulin | Bos taurus | А | 162 | В | 162 |
| 22 | 1BJF | 2.40 | Neurocalcin Delta | Bos taurus | А | 193 | В | 193 |
| 23 | 1BKJ | 1.80 | NADPH-Flavin Oxidoreductase | Vibrio harveyi | А | 240 | В | 240 |

Table 3.3 Dataset of PP complexes from the PIBASE datasetList of PP complexes from the PIBASE dataset were grouped into 335 homodimer complexes (**Table A**).

| 24 | 1BKO | 2.75 | Thymidylate Synthase A | Bacillus subtilis | А | 278 | В | 278 |
|----|------|------|--|--|---|-----|---|-----|
| 25 | 1BQH | 2.80 | Protein (VSV8) | Mus musculus | G | 129 | Н | 129 |
| 26 | 1BUQ | Х | Protein (3-Ketosteroid isomerase-19- nortestosterone-hemisuccinate) | Comamonas testosteroni | А | 125 | В | 125 |
| 27 | 1BVR | 2.80 | Protein (Enoyl-Acyl Carrier Protein (ACP) Reductase) | Mycobacterium tuberculosis | А | 268 | В | 268 |
| 28 | 1BYE | 2.80 | Protein (Glutathione S-transferase) | Zea mays | С | 213 | D | 213 |
| 29 | 1C0T | 2.70 | HIV-1 Reverse Transcriptase | Human immunodeficiency virus 1 | А | 560 | В | 440 |
| 30 | 1CA7 | 2.50 | Protein (Macrophage Migration Inhibitory Factor) | Homo sapiens | А | 114 | С | 114 |
| 31 | 1CBK | 2.00 | Pyrophosphokinase | Haemophilus influenzae | А | 160 | В | 160 |
| 32 | 1CH4 | 2.50 | Module-Substituted Chimera Hemoglobin Beta- Alpha | Homo sapiens | А | 146 | В | 146 |
| 33 | 1CKI | 2.30 | Casein Kinase I Delta | Rattus norvegia | А | 317 | В | 317 |
| 34 | 1CMV | 2.27 | Human Cytomegalovirus Protease | Human herpesvirus 5 | А | 256 | В | 256 |
| 35 | 1CWQ | 2.25 | Bacteriorhodopsin ("M" State Intermediate In Combination With Ground State) | Halobacterium salinarum | А | 248 | В | 248 |
| 36 | 1CYD | 1.80 | Carbonyl Reductase | Mus musculus | А | 244 | В | 244 |
| 37 | 1D0E | 3.00 | Reverse Transcriptase | Moloney murine leukemia virus | А | 259 | В | 259 |
| 38 | 1D0I | 1.80 | Type II 3-Dehydroquinate Hydratase | Streptomyces coelicolor | А | 156 | С | 156 |
| 39 | 1D0O | 1.95 | Bovine Endothelial Nitric Oxide Synthase Heme | Bos taurus | А | 444 | В | 444 |
| 40 | 1D1Z | 1.40 | SAP SH2 Domain | Homo sapiens | А | 104 | В | 104 |
| 41 | 1D4A | 1.70 | Quinone Reductase | Homo sapiens | А | 273 | С | 273 |
| 42 | 1D6S | 2.30 | O-Acetylserine Sulfhydrylase | Salmonella enterica subsp. enterica serovar Typhimurium | А | 322 | В | 322 |
| 43 | 1D6U | 2.40 | Copper Amine Oxidase | Escherichia coli | А | 727 | В | 727 |
| 44 | 1D8U | 2.35 | Non-Symbiotic Hemoglobin | Oryza sativa | А | 166 | В | 166 |
| 45 | 1DBQ | 2.20 | Purine Repressor | Escherichia coli | А | 289 | В | 289 |
| 46 | 1DCL | 2.30 | MCG | Homo sapiens | А | 216 | В | 216 |

| 47 | 1DF8 | 1.51 | Protein (Streptavidin) | Streptomyces avidinii | А | 127 | В | 127 |
|-------|------|------|---|--|---|-----|---|-----|
| 48 | 1DIR | 2.60 | Dihydropteridine Reductase | Rattus norvegia | А | 241 | В | 241 |
| 49 | 1DO5 | 2.75 | Human Copper Chaperone For Superoxide Dismutase Domain | Homo sapiens | А | 154 | С | 154 |
| 50 | 1DOF | 2.10 | Adenylosuccinate Lyase | Pyrobaculum aerophilum | В | 403 | С | 403 |
| 51 | 1DOH | 2.10 | Trihydroxynaphthalene Reductase | Magnaporthe grisea | А | 283 | В | 283 |
| 52 | 1DTY | 2.14 | Adenosylmethionine-8-Amino-7-Oxononanoate Aminotransferase | Escherichia coli | А | 429 | В | 429 |
| 53 | 1DUG | 1.80 | Chimera Of Glutathione S-Transferase-Synthetic Linker-C- Terminal Fibrinogen Gamma Chain | Schistosoma japonicum | А | 234 | В | 234 |
| 54 | 1DVR | 2.36 | Adenylate Kinase | Saccharomyces cerevisiae | А | 220 | В | 220 |
| 55 | 1DZB | 2.00 | SCFV Fragment 1F9 | Mus musculus | А | 253 | В | 253 |
| 56 | 1E3I | 2.08 | Alcohol Dehydrogenase, Class II | Mus musculus | А | 376 | В | 376 |
| 57 | 1E3U | 1.66 | Beta-Lactamase Oxa-10 | Pseudomonas aeruginosa | А | 246 | С | 246 |
| 58 | 1E5E | 2.18 | Methionine Gamma-Lyase | Trichomonas vaginalis G3 | А | 404 | В | 404 |
| 59 | 1E7W | 1.75 | Pteridine Reductase | Leishmania major | А | 291 | В | 291 |
| 60 | 1E8T | 2.50 | Hemagglutinin-Neuraminidase | Newcastle disease virus (strain Kansas) | А | 454 | В | 454 |
| 61 | 1E9S | 2.50 | Conjugal Transfer Protein TRWB | Escherichia coli | Е | 437 | F | 437 |
| 62 | 1ECE | 2.40 | Endocellulase E1 | Acidothermus cellulolyticus | А | 358 | В | 358 |
| 63 | 1ECS | 1.70 | Bleomycin Resistance Protein | Klebsiella pneumoniae | А | 126 | В | 126 |
| 64 | 1EE0 | 2.05 | 2-Pyrone Synthase | Gerbera hybrid cultivar | А | 402 | В | 402 |
| 65 | 1EI1 | 2.30 | DNA Gyrase B | Escherichia coli | А | 391 | В | 391 |
| 66 | 1EK6 | 1.50 | UDP-Galactose 4-Epimerase | Homo sapiens | А | 348 | В | 348 |
| 67 | 1EKX | 1.95 | Aspartate Transcarbamoylase | Escherichia coli | В | 311 | С | 311 |
| 68 | 1EM6 | 2.20 | Liver Glycogen Phosphorylase | Homo sapiens | А | 847 | В | 847 |
| 69 | 1EQU | 3.00 | Protein (Estradiol 17 Beta-Dehydrogenase 1) | Homo sapiens | А | 327 | В | 327 |
| 70 | 1EWK | 2.20 | Metabotropic Glutamate Receptor Subtype 1 | Rattus norvegia | А | 490 | В | 490 |
| ~ · · | 1 | | | | | | | |

| 71 | 1EZR | 2.50 | Nucleoside Hydrolase | Leishmania major | А | 314 | В | 314 |
|--------|------|------|---|-------------------------------|---|-----|---|-----|
| 72 | 1F1J | 2.35 | Caspase-7 Protease | Homo sapiens | А | 305 | В | 305 |
| 73 | 1F2D | 2.00 | 1-Aminocyclopropane-1-Carboxylate Deaminase | Cyberlindnera saturnus | А | 341 | В | 341 |
| 74 | 1F3T | 2.00 | Ornithine Decarboxylase | Trypanosoma brucei | С | 425 | D | 425 |
| 75 | 1F5P | 2.90 | Hemoglobin V | Petromyzon marinus | А | 149 | В | 149 |
| 76 | 1F6B | 1.70 | SAR1 | Cricetulus griseus | А | 198 | В | 198 |
| 77 | 1F6R | 2.20 | Alpha-Lactalbumin | Bos taurus | А | 123 | С | 123 |
| 78 | 1FR1 | 2.00 | Beta-Lactamase | Citrobacter freundii | А | 361 | В | 361 |
| 79 | 1FTM | 1.70 | Glutamate Receptor Subunit 2 | Rattus norvegia | А | 263 | С | 263 |
| 80 | 1FUJ | 2.20 | PR3 | Homo sapiens | А | 221 | D | 221 |
| 81 | 1G1A | 2.50 | DTDP-D-glucose 4,6-Dehydratase Salmonella Enterica | Salmonella enterica | А | 352 | В | 352 |
| 82 | 1G2O | 1.75 | Purine Nucleoside Phosphorylase | Mycobacterium tuberculosis | А | 268 | С | 268 |
| 83 | 1G4A | 3.00 | Atp-Dependent Protease Hslv | Escherichia coli | В | 175 | А | 175 |
| 84 | 1G5Q | 2.57 | Epidermin Modifying Enzyme Epid | Staphylococcus epidermidis | А | 181 | D | 181 |
| 85 | 1G6O | 2.50 | CAG-Alpha | Helicobacter pylori | А | 330 | В | 330 |
| 86 | 1G6Y | 2.80 | URE2 Protein | Saccharomyces cerevisiae | А | 261 | В | 261 |
| 87 | 1G85 | 1.80 | Odorant-Binding Protein | Bos taurus | А | 159 | В | 159 |
| 88 | 1G8Y | 2.40 | Regulatory Protein REPA | Escherichia coli | С | 279 | D | 279 |
| 89 | 1GEG | 1.70 | Acetoin Reductase | Klebsiella pneumoniae | В | 256 | С | 256 |
| 90 | 1GH7 | 3.00 | Cytokine Receptor Common Beta Chain | Homo sapiens | А | 419 | В | 419 |
| 91 | 1GMY | 1.90 | Cathepsin B | Homo sapiens | А | 261 | В | 261 |
| 92 | 1GNX | 1.68 | Beta-Glucosidase | Streptomyces sp. | А | 479 | В | 479 |
| 93 | 1GP7 | 2.60 | Phospholipase A2 | Ophiophagus hannah | А | 151 | В | 151 |
| 94 | 1GPM | 2.20 | GMP Synthetase | Escherichia coli K12 | А | 525 | В | 525 |
| Contin | hore | | | | | | | |

| 95 | 1GR7 | 1.80 | Azurin | Pseudomonas aeruginosa | А | 128 | D | 128 |
|--------|------|------|--|-------------------------------------|---|-----|---|-----|
| 96 | 1GRI | 3.10 | Growth Factor Bound Protein 2 | Homo sapiens | А | 217 | В | 217 |
| 97 | 1GTV | 1.55 | Thymidylate Kinase | Mycobacterium tuberculosis | А | 214 | В | 214 |
| 98 | 1GWN | 2.10 | RHO-Related GTP-Binding Protein Rhoe | Mus musculus | А | 205 | С | 205 |
| 99 | 1GYL | 3.00 | Glycolate Oxidase | Spinacia oleracea | А | 369 | В | 369 |
| 100 | 1GZ6 | 2.38 | Estradiol 17 Beta-Dehydrogenase 4 | Rattus norvegia | А | 319 | В | 319 |
| 101 | 1GZM | 2.65 | Rhodopsin | Bos taurus | А | 349 | В | 349 |
| 102 | 1H48 | 2.30 | 2C-Methyl-D-Erythritol-2,4-Cyclodiphosphate Synthase | Escherichia coli BL21 | А | 161 | В | 161 |
| 103 | 1H4G | 1.10 | Xylanase | (Bacillus) agaradhaerens | А | 207 | В | 207 |
| 104 | 1H5B | 1.85 | Murine T Cell Receptor (TCR) Valpha Domain | Mus musculus | А | 113 | В | 113 |
| 105 | 1H5Q | 1.50 | NADP-Dependent Mannitol Dehydrogenase | Agaricus bisporus | А | 265 | В | 265 |
| 106 | 1H8V | 1.90 | Endo-Beta-1,4-Glucanase | Trichoderma reesei | D | 218 | Е | 218 |
| 107 | 1HAK | 3.00 | Annexin V | Homo sapiens | А | 320 | В | 320 |
| 108 | 1HG4 | 2.40 | Ultraspiracle | Drosophila melanogaster | А | 279 | D | 279 |
| 109 | 1HKV | 2.60 | Diaminopimelate Decarboxylase | Mycobacterium tuberculosis H37RV | А | 453 | В | 453 |
| 110 | 1HP0 | 2.10 | Inosine-Adenosine-Guanosine-Preferring Nucleoside Hydrolase | Trypanosoma vivax | А | 339 | В | 339 |
| 111 | 1HR6 | 2.50 | Mitochondrial Processing Peptidase Alpha Subunit | Saccharomyces cerevisiae | А | 475 | Е | 475 |
| 112 | 1I0Z | 2.10 | L-Lactate Dehydrogenase H Chain | Homo sapiens | А | 333 | В | 333 |
| 113 | 1141 | 3.20 | Cystathionine Gamma-Synthase | Nicotiana tabacum | F | 445 | G | 445 |
| 114 | 1I4U | 1.15 | Crustacyanin | Homarus gammanus | А | 181 | В | 181 |
| 115 | 1171 | 2.35 | Peroxisome Proliferator Activated Receptor Gamma | Homo sapiens | А | 292 | В | 292 |
| 116 | 1185 | 3.20 | T Lymphocyte Activation Antigen CD86 | Homo sapiens | А | 110 | В | 110 |
| 117 | 1116 | 2.10 | Kinesin-Related Motor Protein EG5 | Homo sapiens | А | 368 | В | 368 |
| Contin | ued | | | | | | | |

| 118 | 1IJL | 2.60 | Phospholipase A2 | Deinagkistrodon acutus | А | 123 | В | 123 |
|--------|------|------|--|---|---|-----|---|-----|
| 119 | 1IY8 | 1.60 | Levodione Reductase | Leifsonia aquatica | В | 267 | С | 267 |
| 120 | 1IYK | 2.30 | Myristoyl-COA:Protein N-Myristoyltransferase | Candida albicans | А | 392 | В | 392 |
| 121 | 1IZ1 | 2.50 | LYSR-Type Regulatory Protein | Cupriavidus necator | А | 294 | Р | 294 |
| 122 | 1J7E | 2.55 | Vitamin D Binding Protein | Homo sapiens | А | 458 | В | 458 |
| 123 | 1J90 | 2.56 | Deoxyribonucleoside Kinase | Drosophila melanogaster | А | 230 | В | 230 |
| 124 | 1JDS | 1.80 | 5'-Methylthioadenosine Phosphorylase | Sulfolobus solfatarius | А | 236 | С | 236 |
| 125 | 1JG8 | 1.80 | L-Allo-Threonine Aldolase | Thermotoga maritima | А | 347 | В | 347 |
| 126 | 1JKF | 2.40 | Myo-Inositol-1-Phosphate Synthase | Saccharomyces cerevisiae | А | 533 | В | 533 |
| 127 | 1JSW | 2.70 | L-Aspartate Ammonia-Lyase | Escherichia coli | А | 478 | D | 478 |
| 128 | 1JT0 | 2.90 | Hypothetical Transcriptional Regulator In Qaca 5'Region | Staphylococcus aureus | А | 194 | С | 194 |
| 129 | 1K2O | 1.65 | Cytochrome P450CAM | Pseudomonas putida | А | 414 | В | 414 |
| 130 | 1K3F | 2.50 | Uridine Phosphorylase | Escherichia coli | А | 253 | D | 253 |
| 131 | 1KBU | 2.20 | Cre Recombinase | Enterobacteria phage P1 | А | 349 | В | 349 |
| 132 | 1KCZ | 1.90 | Beta-Methylaspartase | Clostridium tetanomorphum | А | 419 | В | 419 |
| 133 | 1KI9 | 2.76 | Adenylate Kinase | Methanothermococcus thermolithotrophicus | А | 192 | В | 192 |
| 134 | 1KOB | 2.30 | Twitchin | Aplysia californica | А | 387 | В | 387 |
| 135 | 1KRU | 2.80 | Galactoside O-Acetyltransferase | EScherichia coli | А | 203 | С | 203 |
| 136 | 1KXJ | 2.80 | Amidotransferase HISH | Thermotoga maritima | А | 205 | В | 205 |
| 137 | 1L8X | 2.70 | Ferrochelatase | Saccharomyces cerevisiae | А | 362 | В | 362 |
| 138 | 1LBH | 3.20 | Intact Lactose Operon Repressor With Gratuitous Inducer IPTG | EScherichia coli | А | 360 | В | 360 |
| 139 | 1LLU | 2.30 | Alcohol Dehydrogenase | Pseudomonas aeruginosa | Е | 342 | F | 342 |
| 140 | 1LRT | 2.20 | Inosine-5'-Monophosphate Dehydrogenase | Tritrichomonas suis | А | 376 | В | 376 |
| Contin | ued | | | | | | | |

| 141 | 1M43 | 2.40 | Plasmepsin II | Plasmodium falciparum | А | 331 | В | 331 |
|---------|-------|------|---|---------------------------------|---|-----|---|-----|
| 142 | 1M48 | 1.95 | Interleukin-2 | Homo sapiens | А | 133 | В | 133 |
| 143 | 1M7N | 2.70 | Insulin-Like Growth Factor I Receptor | Homo sapiens | А | 322 | В | 322 |
| 144 | 1M7W | 2.80 | Hepatocyte Nuclear Factor 4-Alpha | Rattus rattus | А | 250 | В | 250 |
| 145 | 1MCZ | 2.80 | Benzoylformate Decarboxylase | Pseudomonas putida | А | 582 | В | 582 |
| 146 | 1MPY | 2.80 | Catechol 2,3-Dioxygenase | Pseudomonas putida | А | 307 | С | 307 |
| 147 | 1MRU | 3.00 | Probable Serine/Threonine-Protein Kinase Pknb | Mycobacterium tuberculosis | А | 311 | В | 311 |
| 148 | 1MZJ | 2.10 | Beta-Ketoacylsynthase III | Streptomyces sp. R1128 | А | 339 | В | 339 |
| 149 | 1MZN | 1.90 | RXR Retinoid X Receptor | Homo sapiens | А | 240 | С | 240 |
| 150 | 1N0H | 2.80 | Acetolactate Synthase | Saccharomyces cerevisiae | А | 677 | В | 677 |
| 151 | 1N0S | 2.00 | Bilin-Binding Protein | Pieris brassicae | А | 184 | В | 184 |
| 152 | 1N1A | 2.40 | FKBP52 | Homo sapiens | А | 140 | В | 140 |
| 153 | 1N4O | 1.85 | L2 beta-lactamase | Stenotrophomonas maltophilia | А | 276 | В | 276 |
| 154 | 1N7G | 2.20 | GDP-D-Mannose-4,6-Dehydratase | Homo sapiens | А | 381 | В | 381 |
| 155 | 1N9E | 1.65 | Lysyl Oxidase | Komagataella pastoris | А | 787 | В | 787 |
| 156 | 1NBU | 1.60 | Probable Dihydroneopterin Aldolase | Mycobacterium tuberculosis | В | 119 | С | 119 |
| 157 | 1NFQ | 2.40 | Putative Oxidoreductase Rv2002 | Mycobacterium tuberculosis | А | 260 | В | 260 |
| 158 | 1NKS | 2.57 | Adenylate Kinase | Sulfolobus acidocaldarius | А | 194 | В | 194 |
| 159 | 1NTO | 1.94 | NAD-Dependent Alcohol Dehydrogenase | Sulfolobus solfataricus | С | 347 | D | 347 |
| 160 | 1NW4 | 2.20 | Uridine Phosphorylase, Putative | Plasmodium falciparum 3D7 | А | 276 | В | 276 |
| 161 | 1NX9 | 2.20 | Alpha-Amino Acid Ester Hydrolase | Acetobacter pasteurianus | В | 652 | D | 652 |
| 162 | 104S | 1.90 | Aspartate Aminotransferase | Thermotoga maritima | А | 389 | В | 389 |
| 163 | 1O5I | 2.50 | 3-Oxoacyl-(Acyl Carrier Protein) Reductase | Thermotoga maritima | А | 249 | D | 249 |
| 7 antin | h a a | | | | | | | |

| 164 | 1050 | 2.30 | Uracil Phosphoribosyltransferase | Thermotoga maritima | А | 221 | В | 221 |
|-----|------|------|---|--|---|-----|---|-----|
| 165 | 1061 | 1.90 | Aminotransferase | Campylobacter jejuni | А | 394 | В | 394 |
| 166 | 1063 | 2.00 | ATP Phosphoribosyltransferase | THermotoga maritima | А | 219 | В | 219 |
| 167 | 106E | 2.30 | Capsid Protein P40 | Human herpesvirus 4 | А | 235 | В | 235 |
| 168 | 109J | 2.40 | Aldehyde Dehydrogenase, Cytosolic 1 | Elephantulus edwardii | А | 501 | В | 501 |
| 169 | 10AT | 2.50 | Ornithine Aminotransferase | Homo sapiens | В | 439 | С | 439 |
| 170 | 10D2 | 2.70 | Acetyl-Coenzyme A Carboxylase | Saccharomyces cerevisiae | А | 805 | В | 805 |
| 171 | 10DL | 2.10 | Purine Nucleoside Phosphorylase | Thermus thermophilus HB8 | А | 235 | В | 235 |
| 172 | 10E7 | 1.80 | Glutathione S-Transferase | Schistosoma haematobium | А | 211 | В | 211 |
| 173 | 1OJ6 | 1.95 | Neuroglobin | Homo sapiens | А | 151 | В | 151 |
| 174 | 10ME | 2.30 | Beta-Lactamase | Staphylococcus aureus | А | 258 | В | 258 |
| 175 | 10PL | 3.42 | Proto-Oncogene Tyrosine-Protein Kinase | Homo sapiens | А | 537 | В | 537 |
| 176 | 10RR | 1.50 | CDP-Tyvelose-2-Epimerase | Salmonella enterica subsp. enterica serovar Typhi | А | 347 | В | 347 |
| 177 | 10RT | 3.00 | Ornithine Transcarbamoylase | Pseudomonas aeruginosa | С | 335 | D | 335 |
| 178 | 10RW | 2.84 | Dipeptidyl Peptidase IV | Sus scrofa | В | 728 | D | 728 |
| 179 | 10VL | 2.20 | Orphan Nuclear Receptor NURR1 (MSE 414, 496, 511) | Homo sapiens | А | 271 | В | 271 |
| 180 | 10VM | 2.65 | Indole-3-Pyruvate Decarboxylase | Enterobacter cloacae | В | 552 | D | 552 |
| 181 | 10YJ | 1.95 | Glutathione S-Transferase | Oryza sativa | С | 231 | D | 231 |
| 182 | 10ZF | 2.30 | Acetolactate Synthase, Catabolic | Klebsiella pneumoniae | А | 566 | В | 566 |
| 183 | 1P0K | 1.90 | Isopentenyl-Diphosphate Delta-Isomerase | Bacillus subtilis | А | 349 | В | 349 |
| 184 | 1P4E | 2.70 | Recombinase FLP Protein | Saccharomyces cerevisiae | А | 429 | D | 429 |
| 185 | 1P60 | 1.96 | Deoxycytidine Kinase | Homo sapiens | А | 263 | В | 263 |
| 186 | 1P7C | 2.10 | Thymidine Kinase | Herpes simplex virus (type 1 / strain 17) | А | 343 | В | 343 |

| 187 | 1P93 | 2.70 | Glucocorticoid Receptor | Homo sapiens | А | 280 | С | 280 |
|------|------|------|---|--|---|-----|---|-----|
| 188 | 1PKG | 2.90 | C-Kit Protein | Homo sapiens | А | 329 | В | 329 |
| 189 | 1PV4 | 3.00 | Transcription Termination Factor Rho | EScherichia coli | С | 419 | D | 419 |
| 190 | 1PWE | 2.80 | L-Serine Dehydratase | Rattus norvegia | С | 327 | D | 327 |
| 191 | 1PWX | 1.80 | Halohydrin Dehalogenase | Agrobacterium tumefaciens | В | 254 | С | 254 |
| 192 | 1PY4 | 2.90 | Beta-2-Microglobulin Precursor | Homo sapiens | А | 100 | С | 100 |
| 193 | 1Q0C | 2.10 | Homoprotocatechuate 2,3-Dioxygenase | Brevibacterium fuscum | А | 365 | С | 365 |
| 194 | 1Q3D | 2.20 | Glycogen Synthase Kinase-3 Beta | Homo sapiens | А | 424 | В | 424 |
| 195 | 1Q3G | 2.65 | UDP-N-Acetylglucosamine 1- Carboxyvinyltransferase | Enterobacter cloacae | А | 419 | В | 419 |
| 196 | 1Q57 | 3.45 | DNA Primase/Helicase | Enterobacteria phage T7 | А | 503 | С | 503 |
| 197 | 1Q5Q | 2.60 | Proteasome Alpha-Type Subunit 1 | Rhodococcus erythropolis | С | 259 | В | 259 |
| 198 | 1Q6T | 2.30 | Protein-Tyrosine Phosphatase, Non-Receptor Type 1 | Homo sapiens | А | 310 | В | 310 |
| 199 | 1Q8M | 2.60 | Triggering Receptor Expressed On Myeloid Cells 1 | Homo sapiens | А | 127 | В | 127 |
| 200 | 1QPB | 2.40 | Pyruvate Decarboxylase (Form B) | Saccharomyces pastorianus Weihenstephan 34/70 | А | 563 | В | 563 |
| 201 | 1QZF | 2.80 | Bifunctional Dihydrofolate Reductase-Thymidylate Synthase | Cryptosporidium hominis | С | 521 | D | 521 |
| 202 | 1R5K | 2.70 | Estrogen Receptor | Homo sapiens | А | 261 | С | 261 |
| 203 | 1RD5 | 2.02 | Tryptophan Synthase Alpha Chain, Chloroplast | Zea mays | А | 262 | В | 262 |
| 204 | 1RD7 | 2.60 | Dihydrofolate Reductase | EScherichia coli | А | 159 | В | 159 |
| 205 | 1RE5 | 2.60 | 3-Carboxy-Cis,Cis-Muconate Cycloisomerase | Pseudomonas putida KT2440 | В | 450 | С | 450 |
| 206 | 1RJN | 2.30 | menB | Mycobacterium tuberculosis | А | 339 | В | 339 |
| 207 | 1RKX | 1.80 | CDP-Glucose-4,6-Dehydratase | Yersinia pseudotuberculosis | А | 357 | D | 357 |
| 208 | 1RPY | 2.30 | Adaptor Protein APS | Rattus norvegia | А | 114 | В | 114 |
| 209 | 1RQR | 2.67 | 5'-Fluoro-5'-Deoxyadenosine Synthase | Streptomyces cattleya | А | 299 | С | 299 |
| a .• | 1 | | | | | | | |

| 210 | 1S2G | 2.10 | Purine Trans Deoxyribosylase | Lactobacillus helveticus | А | 167 | В | 167 |
|-----|------|------|---|--------------------------|---|-----|---|-----|
| 211 | 1SF2 | 2.40 | 4-Aminobutyrate Aminotransferase | EScherichia coli | А | 426 | В | 426 |
| 212 | 1SGF | 3.15 | Nerve Growth Factor | Mus musculus | G | 237 | Z | 237 |
| 213 | 1SL6 | 2.25 | C-Type Lectin DC-Signr | Homo sapiens | D | 184 | А | 184 |
| 214 | 1SN0 | 1.90 | Transthyretin | Sparus aurata | А | 130 | С | 130 |
| 215 | 1SP8 | 2.00 | 4-Hydroxyphenylpyruvate Dioxygenase | Zea mays | А | 418 | В | 418 |
| 216 | 1SPI | 2.80 | Fructose 1,6-Bisphosphatase | Spinacia oleracea | А | 358 | В | 358 |
| 217 | 1SQI | 2.15 | 4-Hydroxyphenylpyruvic Acid Dioxygenase | Rattus norvegia | А | 393 | В | 393 |
| 218 | 1SQL | 2.20 | Dihydroneopterin Aldolase | Arabidopsis thaliana | D | 146 | Е | 146 |
| 219 | 1SU2 | 2.70 | Mutt/Nudix Family Protein | Deinococcus radiodurans | А | 159 | В | 159 |
| 220 | 1SXG | 2.75 | Glucose-Resistance Amylase Regulator | Bacillus megaterium | А | 280 | D | 280 |
| 221 | 1T2O | 2.30 | Sortase | Staphylococcus aureus | В | 146 | А | 146 |
| 222 | 1T8S | 2.60 | Amp Nucleosidase | EScherichia coli | А | 484 | В | 484 |
| 223 | 1T91 | 1.90 | Ras-Related Protein Rab-7 | Homo sapiens | А | 207 | D | 207 |
| 224 | 1T97 | 2.70 | Lysozyme | Enterobacteria phage T4 | А | 175 | В | 175 |
| 225 | 1TAH | 3.00 | Lipase | Burkholderia glumae | А | 318 | В | 318 |
| 226 | 1TC0 | 2.20 | Endoplasmin | Canis lupus familiaris | А | 236 | В | 236 |
| 227 | 1TEX | 2.60 | Stf0 Sulfotransferase | Mycobacterium smegmatis | А | 287 | В | 287 |
| 228 | 1TF7 | 2.80 | KaiC | Synechococcus sp. | В | 525 | С | 525 |
| 229 | 1TFC | 2.40 | Estrogen-Related Receptor Gamma | Homo sapiens | А | 251 | В | 251 |
| 230 | 1TPZ | 2.00 | Interferon-Inducible GTPase | Mus musculus | А | 422 | В | 422 |
| 231 | 1TW2 | 2.50 | Carminomycin 4-O-Methyltransferase | Streptomyces peucetius | А | 360 | В | 360 |
| 232 | 1TXT | 2.50 | 3-Hydroxy-3-Methylglutaryl-CoA Synthase | Staphylococcus aureus | А | 388 | В | 388 |
| a | 1 | | | | | | | |

| 233 | 1TZ3 | 2.90 | Putative Sugar Kinase | Salmonella enterica subsp. enterica serovar Typhimurium str. LT2 | А | 339 | В | 339 |
|--------|------|------|--|---|---|-----|---|-----|
| 234 | 1U0L | 2.80 | Probable GTPase engC | Thermotoga maritima | А | 301 | С | 301 |
| 235 | 1U1I | 1.90 | Myo-Inositol-1-Phosphate Synthase | Archaeoglobus fulgidus DSM 4304 | А | 392 | В | 392 |
| 236 | 1U46 | 2.00 | Activated CDC42 Kinase 1 | Homo sapiens | А | 291 | В | 291 |
| 237 | 1U5Q | 2.10 | Serine/Threonine Protein Kinase TAO2 | Rattus norvegia | А | 348 | В | 348 |
| 238 | 1UB7 | 2.30 | 3-Oxoacyl-[Acyl-Carrier Protein] Synthase | Thermus thermophilus HB8 | А | 322 | В | 322 |
| 239 | 1UFH | 2.20 | YYCN Protein | Bacillus subtilis subsp. subtilis str. 168 | В | 180 | А | 180 |
| 240 | 1UI5 | 2.40 | A-Factor Receptor Homolog | Streptomyces coelicolor A3(2) | А | 215 | В | 215 |
| 241 | 1UIM | 2.15 | Threonine Synthase | Thermus thermophilus | А | 351 | В | 351 |
| 242 | 1UIU | 1.85 | Nickel-Binding Periplasmic Protein | EScherichia coli | А | 502 | В | 502 |
| 243 | 1UMO | 2.59 | Cytoglobin | Homo sapiens | А | 190 | В | 190 |
| 244 | 1UPA | 2.35 | Carboxyethylarginine Synthase | Streptomyces clavuligerus | В | 573 | D | 573 |
| 245 | 1URZ | 2.70 | Envelope Protein | Tick-borne encephalitis virus (WESTERN SUBTYPE) | В | 401 | С | 401 |
| 246 | 1USI | 1.80 | Leucine-Specific Binding Protein | Escherichia coli | А | 346 | С | 346 |
| 247 | 1UTR | Х | Uteroglobin | Rattus norvegia | А | 96 | В | 96 |
| 248 | 1UU0 | 2.85 | Histidinol-Phosphate Aminotransferase | Thermotoga maritima | А | 335 | В | 335 |
| 249 | 1UZM | 1.49 | 3-Oxoacyl-[Acyl-Carrier Protein] Reductase | Mycobacterium tuberculosis H37Rv | А | 247 | В | 247 |
| 250 | 1V2I | 2.20 | Hemagglutinin-Neuraminidase Glycoprotein | Human parainfluenza virus 3 | А | 431 | В | 431 |
| 251 | 1V4J | 2.85 | Octoprenyl-Diphosphate Synthase | Thermotoga maritima | А | 299 | В | 299 |
| 252 | 1VC8 | 2.00 | NDX1 | Thermus thermophilus HB8 | А | 126 | В | 126 |
| 253 | 1VDM | 2.50 | Purine Phosphoribosyltransferase | Pyrococcus horikoshii | В | 153 | Н | 153 |
| 254 | 1VDW | 1.30 | Hypothetical Protein PH1897 | Pyrococcus horikoshii | А | 254 | В | 254 |
| 255 | 1VEA | 2.80 | Hut Operon Positive Regulatory Protein | Bacillus subtilis | А | 148 | В | 148 |
| Contin | ued | | | | | | | |

| 1VIQ 1VIY 1VKG 1VM6 1VPX 1W0M 1W2Z 1W59 | 2.40 1.89 2.20 2.27 2.40 2.50 2.24 | ADP-Ribose Pyrophosphatase Dephospho-CoA Kinase Histone Deacetylase 8 Dihydrodipicolinate Reductase Protein (Transaldolase (EC 2.2.1.2)) Triosephosphate Isomerase Amine Oxidase, Copper Containing | Escherichia coli Escherichia coli Homo sapiens THermotoga maritima THermotoga maritima Thermoproteus tenax Pisum satiuum | B A B B B | 220 218 377 228 230 226 | C B D C D | 220 218 377 228 230 |
|--|--|---|--|--|--|--|---|
| 1VIY 1VKG 1VM6 1VPX 1W0M 1W2Z 1W59 | 1.89 2.20 2.27 2.40 2.50 2.24 | Dephospho-CoA Kinase Histone Deacetylase 8 Dihydrodipicolinate Reductase Protein (Transaldolase (EC 2.2.1.2)) Triosephosphate Isomerase Amine Oxidase, Copper Containing | Escherichia coli Homo sapiens THermotoga maritima THermotoga maritima Thermoproteus tenax Pisum satiuum | A A B B B | 218 377 228 230 226 | B D C D | 218 377 228 230 |
| 1VKG 1VM6 1VPX 1W0M 1W2Z 1W59 | 2.20 2.27 2.40 2.50 2.24 | Histone Deacetylase 8 Dihydrodipicolinate Reductase Protein (Transaldolase (EC 2.2.1.2)) Triosephosphate Isomerase Amine Oxidase, Copper Containing | Homo sapiens THermotoga maritima THermotoga maritima Thermoproteus tenax Pisum satiuum | A B B B | 377 228 230 226 | B D C D | 377 228 230 |
| 1VM6 1VPX 1W0M 1W2Z 1W59 | 2.27 2.40 2.50 2.24 | Dihydrodipicolinate Reductase Protein (Transaldolase (EC 2.2.1.2)) Triosephosphate Isomerase Amine Oxidase, Copper Containing | THermotoga maritima THermotoga maritima Thermoproteus tenax Pisum satiuum | B B B | 228 230 226 | D C D | 228 230 |
| 1VPX 1W0M 1W2Z 1W59 | 2.40 2.50 2.24 | Protein (Transaldolase (EC 2.2.1.2)) Triosephosphate Isomerase Amine Oxidase, Copper Containing | THermotoga maritima Thermoproteus tenax Pisum satiuum | B B | 230 226 | C D | 230 |
| 1W0M 1W2Z 1W59 | 2.50 2.24 | Triosephosphate Isomerase Amine Oxidase, Copper Containing | Thermoproteus tenax | В | 226 | D | |
| 1W2Z 1W59 | 2.24 | Amine Oxidase, Copper Containing | Disum sativum | | | - | 226 |
| 1W59 | 0.70 | | Fisuin sauvum | A | 649 | В | 649 |
| | 2.70 | Cell Division Protein FTSZ Homolog 1 | Methanocaldococcus jannaschii | А | 364 | В | 364 |
| 1W6T | 2.10 | Enolase | Streptococcus pneumoniae TIGR4 | А | 444 | В | 444 |
| 1W6U | 1.75 | 2,4-Dienoyl-Coa Reductase, Mitochondrial Precursor | Homo sapiens | А | 302 | В | 302 |
| 1W91 | 2.20 | Beta-Xylosidase | Geobacillus stearothermophilus | F | 503 | G | 503 |
| 1WC1 | 1.93 | Adenylate Cyclase | Arthrospira platensis | В | 226 | С | 226 |
| 1WUE | 2.10 | Mandelate Racemase/Muconate Lactonizing Enzyme Family Protein | Enterococcus faecalis V583 | А | 386 | В | 386 |
| 1WYE | 2.80 | 2-Keto-3-Deoxygluconate Kinase | Sulfolobus tokodaii str. 7 | А | 311 | В | 311 |
| 1WZ8 | 1.80 | Enoyl-CoA Hydratase | Thermus thermophilus HB8 | А | 264 | В | 264 |
| 1X27 | 2.70 | Proto-Oncogene Tyrosine-Protein Kinase LCK | Homo sapiens | А | 167 | Е | 167 |
| 1XCA | 2.30 | Cellular Retinoic Acid Binding Protein Type II | Homo sapiens | А | 137 | В | 137 |
| 1XED | 1.90 | Polymeric-Immunoglobulin Receptor | Homo sapiens | А | 117 | В | 117 |
| 1XG5 | 1.53 | ARPG836 | Homo sapiens | А | 279 | В | 279 |
| 1XGM | 2.80 | Methionine Aminopeptidase | Pyrococcus furiosus | А | 295 | В | 295 |
| 1XI9 | 2.33 | Putative Transaminase | Pyrococcus furiosus DSM 3638 | А | 406 | В | 406 |
| 1XKQ | 2.10 | Short-Chain Reductase Family Member (5D234) | Caenorhabditis elegans | А | 280 | В | 280 |
| | 1W59 1W6T 1W6U 1W91 1WC1 1WUE 1WVE 1WZ8 1X27 1XCA 1XED 1XG5 1XGM 1XI9 1XKQ | 1W592.701W6T2.101W6U1.751W912.201WC11.931WUE2.101WYE2.801WZ81.801X272.701XCA2.301XED1.901XG51.531XGM2.801X192.331XKQ2.10 | W222.24Annue Oxidase, Copper Containing1W592.70Cell Division Protein FTSZ Homolog 11W6T2.10Enolase1W6U1.752,4-Dienoyl-Coa Reductase, Mitochondrial Precursor1W912.20Beta-Xylosidase1WC11.93Adenylate Cyclase1WUE2.10Mandelate Racemase/Muconate Lactonizing Enzyme Family Protein1WYE2.802-Keto-3-Deoxygluconate Kinase1WZ81.80Enoyl-CoA Hydratase1X272.70Proto-Oncogene Tyrosine-Protein Kinase LCK1XCA2.30Cellular Retinoic Acid Binding Protein Type II1XED1.90Polymeric-Immunoglobulin Receptor1XG51.53ARPG8361XGM2.80Methionine Aminopeptidase1XI92.33Putative Transaminase1XKQ2.10Short-Chain Reductase Family Member (5D234) | 1W2Z2.24Amine Oxidase, Copper ContainingPisum sativum Methanocaldococcus jannaschii1W592.70Cell Division Protein FTSZ Homolog 1Methanocaldococcus jannaschii1W6T2.10EnolaseStreptococcus pneumoniae TIGR41W6U1.752,4-Dienoyl-Coa Reductase, Mitochondrial PrecursorHomo sapiens Geobacillus stearothermophilus1W912.20Beta-XylosidaseGeobacillus stearothermophilus1WC11.93Adenylate CyclaseArthrospira platensis1WUE2.10Mandelate Racemase/Muconate Lactonizing Enzyme Family ProteinEnterococcus faecalis V5831WYE2.802-Keto-3-Deoxygluconate KinaseSulfolobus tokodaii str. 71WZ81.80Enoyl-CoA HydrataseHB81X272.70Proto-Oncogene Tyrosine-Protein Kinase LCKHomo sapiens1XED1.90Polymeric-Immunoglobulin ReceptorHomo sapiens1XG51.53ARPG836Homo sapiens1XGM2.80Methionine AminopeptidasePyrococcus furiosus AG381XI92.33Putative TransaminasePyrococcus furiosus DSM 36381XKQ2.10Short-Chain Reductase Family Member (5D234)Caenorhabditis elegans | 1W2Z2.24Amine Oxidase, Copper ContainingPisum sativumA1W592.70Cell Division Protein FTSZ Homolog 1Methanocaldococcus jannaschiiA1W6T2.10EnolaseStreptococcus pneumoniae TIGR4A1W6U1.752,4-Dienoyl-Coa Reductase, Mitochondrial PrecursorHomo sapiensA1W912.20Beta-XylosidaseGeobacillus stearothermophilusF1WC11.93Adenylate CyclaseArthrospira platensisB1WUE2.10Mandelate Racemase/Muconate Lactonizing Enzyme Family ProteinFnerococcus faecalis V583A1WYE2.802-Keto-3-Deoxygluconate KinaseSulfolobus tokodaii str. 7 HB8A1WZ81.80Enoyl-CoA HydrataseThermus thermophilus HB8A1X272.70Proto-Oncogene Tyrosine-Protein Kinase LCKHomo sapiensA1XED1.90Polymeric-Immunoglobulin ReceptorHomo sapiensA1XG51.53ARPG836Homo sapiensA1XG82.80Methionine AminopeptidasePyrococcus furiosusA1XI92.33Putative TransaminasePyrococcus furiosus DSM 3638A1XKQ2.10Short-Chain Reductase Family Member (5D234)Caenorhabditis elegansA | 11W2Z2.24Amine Oxidase, Copper ContainingPisum sativumA6491W592.70Cell Division Protein FTSZ Homolog 1Methanocaldococcus jannaschiiA3641W6T2.10EnolaseStreptococcus pneumoniae TIGR4A4441W6U1.752,4-Dienoyl-Coa Reductase, Mitochondrial PrecursorHomo sapiensA3021W912.20Beta-XylosidaseGeobacillus stearothermophilusF5031WC11.93Adenylate CyclaseArthrospira platensisB2261WUE2.10Mandelate Racemase/Muconate Lactonizing Enzyme Family ProteinEnterococcus faecalis V583A3861WYE2.802-Keto-3-Deoxygluconate KinaseSulfolobus tokodaii str. 7A3111WZ81.80Enoyl-CoA HydrataseThermus thermophilus HB8A2641X272.70Proto-Oncogene Tyrosine-Protein Kinase LCKHomo sapiensA1671XCA2.30Cellular Retinoic Acid Binding Protein Type IIHomo sapiensA1371XED1.90Polymeric-Immunoglobulin ReceptorHomo sapiensA1171XG51.53ARPG836Homo sapiensA2951XI92.33Putative TransaminasePyrococcus furiosus DSM 3638A4061XKQ2.10Short-Chain Reductase Family Member (5D234)Caenorhabditis elegansA280 | TwoM2.50Trosephosphate IsomeraseThermoproteus tenaxB226D1W2Z2.24Amine Oxidase, Copper ContainingPisum sativumA649B1W592.70Cell Division Protein FTSZ Homolog 1Methanocaldococcus jannaschiiA364B1W6T2.10EnolaseStreptococcus pneumoniae TIGR4A444B1W6U1.752,4-Dienoyl-Coa Reductase, Mitochondrial PrecursorHomo sapiensA302B1W912.20Beta-XylosidaseGeobacillus stearothermophilusF503G1WC11.93Adenylate CyclaseArthrospira platensisB226C1WUE2.10Mandelate Racemase/Muconate Lactonizing Enzyme Family ProteinEnterococcus faecalis V583A386B1WZ81.80Enoyl-CoA HydrataseSulfolobus tokodaii str. 7A311B1WZ81.80Enoyl-CoA HydrataseHomo sapiensA167E1XCA2.30Cellular Retinoic Acid Binding Protein Type IIHomo sapiensA137B1XED1.90Polymeric-Immunoglobulin ReceptorHomo sapiensA117B1XGS1.53ARPG836Homo sapiensA295B1XI92.33Putative TransaminasePyrococcus furiosus DSM 3638A260B1XKQ2.10Short-Chain Reductase Family Member (5D234)Caenorhabditis elegansA280B |

| 279 | 1XNX | 2.90 | Constitutive Androstane Receptor | Mus musculus | А | 256 | В | 256 |
|--------|------|------|--|--------------------------------|---|------|---|------|
| 280 | 1XSE | 2.90 | 11beta-Hydroxysteroid Dehydrogenase Type 1 | Cavia porcellus | А | 295 | В | 295 |
| 281 | 1XTT | 1.80 | Probable Uracil Phosphoribosyltransferase | Sulfolobus solfataricus | В | 216 | С | 216 |
| 282 | 1XYG | 2.19 | Putative N-Acetyl-Gamma-Glutamyl-Phosphate Reductase | Arabidopsis thaliana | В | 359 | С | 359 |
| 283 | 1Y1M | 1.80 | Glutamate [NMDA] Receptor Subunit Zeta 1 | Rattus norvegia | А | 292 | В | 292 |
| 284 | 1Y1P | 1.60 | Aldehyde Reductase II | Sporidiobolus salmonicolor | А | 342 | В | 342 |
| 285 | 1YDE | 2.40 | Retinal Dehydrogenase/Reductase 3 | Homo sapiens | В | 270 | Р | 270 |
| 286 | 1YO6 | 2.60 | Putative Carbonyl Reductase Sniffer | Caenorhabditis elegans | А | 250 | В | 250 |
| 287 | 1YP2 | 2.11 | Glucose-1-Phosphate Adenylyltransferase Small Subunit | Solanum tuberosum | С | 451 | D | 451 |
| 288 | 1YQ2 | 1.90 | Beta-Galactosidase | Arthrobacter sp. C2-2 | А | 1024 | В | 1024 |
| 289 | 1YTA | 2.20 | Oligoribonuclease | Escherichia coli | А | 180 | В | 180 |
| 290 | 1YXM | 1.90 | Peroxisomal Trans 2-Enoyl Coa Reductase | Homo sapiens | А | 303 | В | 303 |
| 291 | 1Z08 | 1.80 | Ras-Related Protein Rab-21 | Homo sapiens | В | 170 | D | 170 |
| 292 | 1Z1B | 3.80 | Integrase | Enterobacteria phage lambda | А | 356 | В | 356 |
| 293 | 1Z5A | 2.20 | Type II DNA Topoisomerase VI Subunit B | Sulfolobus shibatae | А | 469 | В | 469 |
| 294 | 1ZEM | 1.90 | Xylitol Dehydrogenase | Gluconobacter oxydans | Е | 262 | G | 262 |
| 295 | 1ZPD | 1.86 | Pyruvate Decarboxylase | Zymomonas mobilis | В | 568 | E | 568 |
| 296 | 2A2G | 2.90 | Protein (Alpha-2u-Globulin) | Rattus norvegia | А | 181 | С | 181 |
| 297 | 2A7K | 2.24 | CarB | Pectobacterium carotovorum | G | 250 | Н | 250 |
| 298 | 2AAW | 2.40 | Glutathione S-Transferase | Plasmodium falciparum | А | 222 | С | 222 |
| 299 | 2AG5 | 1.84 | Dehydrogenase/Reductase (SDR Family) Member 6 | Homo sapiens | А | 246 | В | 246 |
| 300 | 2ANC | 3.20 | Guanylate Kinase | Escherichia coli | С | 207 | Е | 207 |
| 301 | 2B4G | 1.95 | Dihydroorotate Dehydrogenase | Trypanosoma brucei | А | 317 | В | 317 |
| Contin | hed | | | | | | | |

| 302 | 2BCK | 2.80 | HLA Class I Histocompatibility Antigen, A-24 Alpha | Homo sapiens | А | 294 | D | 294 |
|-----|------|------|--|---|---|-----|---|-----|
| 303 | 2BD0 | 1.70 | Sepiapterin Reductase | Chlorobium tepidum TLS | А | 244 | В | 244 |
| 304 | 2BTW | 2.00 | ALR0975 Protein | Nostoc sp. PCC 7120 | A | 254 | B | 254 |
| 305 | 2BVC | 2.10 | Glutamine Synthetase 1 | Mycobacterium tuberculosis H37Rv | D | 486 | E | 486 |
| 306 | 2C7F | 2.70 | Alpha-L-Arabinofuranosidase | Ruminiclostridium thermocellum | А | 513 | D | 513 |
| 307 | 2CLT | 2.67 | Interstitial Collagenase | Homo sapiens | А | 367 | В | 367 |
| 308 | 2D1Y | 1.65 | Hypothetical Protein TT0321 | Thermus thermophilus | В | 256 | С | 256 |
| 309 | 2DFT | 2.80 | Shikimate Kinase | Mycobacterium tuberculosis | А | 176 | С | 176 |
| 310 | 2DNS | 2.40 | D-Amino Acid Amidase | Ochrobactrum anthropi | А | 363 | В | 363 |
| 311 | 2DW6 | 2.30 | BII6730 Protein | Bradyrhizobium japonicum | А | 389 | В | 389 |
| 312 | 2DY4 | 2.65 | DNA Polymerase | Enterobacteria phage RB69 | D | 903 | В | 903 |
| 313 | 2EGH | 2.20 | 1-Deoxy-D-Xylulose 5-Phosphate Reductoisomerase | Escherichia coli str. K-12 substr. W3110 | В | 424 | А | 424 |
| 314 | 2EVO | 1.70 | Methionine Aminopeptidase | Escherichia coli | А | 264 | В | 264 |
| 315 | 2EW8 | 2.10 | (S)-1-Phenylethanol Dehydrogenase | Aromatoleum aromaticum EbN1 | В | 249 | D | 249 |
| 316 | 2F1G | 1.90 | Cathepsin S | Homo sapiens | А | 220 | В | 220 |
| 317 | 2F73 | 2.50 | Fatty Acid-Binding Protein, Liver | Homo sapiens | А | 149 | В | 149 |
| 318 | 2FG6 | 2.80 | Putative Ornithine Carbamoyltransferase | Bacteroides fragilis NCTC 9343 | С | 338 | D | 338 |
| 319 | 2FM6 | 1.75 | Metallo-Beta-Lactamase L1 | Stenotrophomonas maltophilia | А | 269 | В | 269 |
| 320 | 2FYI | 2.80 | Hth-Type Transcriptional Regulator CBL | Escherichia coli K-12 | А | 228 | В | 228 |
| 321 | 2GF0 | 1.90 | GTP-Binding Protein DI-RAS1 | Homo sapiens | А | 199 | В | 199 |
| 322 | 2GIC | 2.92 | Nucleocapsid Protein | Vesicular stomatitis Indiana virus | А | 422 | В | 422 |
| 323 | 216A | 2.20 | Adenosine Kinase | Homo sapiens | А | 345 | D | 345 |
| 324 | 2IFA | 2.30 | Hypothetical Protein SMU.260 | Streptococcus mutans UA159 | А | 208 | В | 208 |
| ~ | - | | | | | | | |

| 325 | 2IGO | 1.95 | Pyranose Oxidase | Trametes ochracea | С | 623 | D | 623 |
|-----|------|------|---|------------------------------|---|-----|---|-----|
| 326 | 2IJ2 | 1.20 | Cytochrome P450 BM3 | Bacillus megaterium | В | 470 | А | 470 |
| 327 | 2NUU | 2.50 | Ammonia Channel | Escherichia coli | В | 415 | С | 415 |
| 328 | 2023 | 1.20 | HADH2 Protein | Homo sapiens | А | 265 | В | 265 |
| 329 | 202Y | 2.20 | Enoyl-Acyl Carrier Reductase | Plasmodium falciparum 3D7 | С | 349 | D | 349 |
| 330 | 20KR | 2.00 | Mitogen-Activated Protein Kinase 14 | Homo sapiens | А | 366 | D | 366 |
| 331 | 2PD3 | 2.50 | Enoyl-[Acyl-Carrier-Protein] Reductase [NADH] | Helicobacter pylori | В | 275 | С | 275 |
| 332 | 2PMT | 2.70 | Glutathione Transferase | Proteus mirabilis | С | 203 | D | 203 |
| 333 | 2TAA | 3.00 | TAKA-Amylase A | Aspergillus oryzae | А | 478 | В | 478 |
| 334 | 3TAT | 3.50 | Tyrosine Aminotransferase | Escherichia coli | С | 397 | D | 397 |
| 335 | 7MDH | 2.40 | Protein (Malate Dehydrogenase) | Sorghum bicolor | D | 375 | С | 375 |

| | | | Resolution | Chain | | | Chain | | |
|---|-------|--------|------------|-------|---|--------|-------|---|--------|
| _ | No | PDB ID | (Å) | one | Name of chain one | Length | two | Name of chain two | Length |
| | 1 | 1A22 | 2.60 | В | Growth Hormone Receptor | 238 | А | Growth Hormone | 191 |
| | 2 | 1A2K | 2.50 | С | RAN | 216 | В | Nuclear Transport Factor 2 | 127 |
| | 3 | 1A5A | 1.90 | А | Tryptophan Synthase Alpha chain | 268 | В | Tryptophan Synthase Beta Chain | 397 |
| | 4 | 1A5A | 1.90 | В | Tryptophan Synthase Beta chain | 397 | А | Tryptophan Synthase Alpha Chain | 268 |
| | 5 | 1ACB | 2.00 | E | Alpha-Chymotrypsin | 241 | I | Eglin C | 63 |
| | 6 | 1AIG | 2.60 | Ρ | Photosynthetic Reaction Center (H Subunit) | 260 | 0 | Photosynthetic Reaction Center (M Subunit) | 307 |
| | 7 | 1AIG | 2.60 | L | Photosynthetic Reaction Center (L Subunit) | 281 | н | Photosynthetic Reaction Center (H Subunit) | 260 |
| | 8 | 1AIP | 3.00 | А | Elongation Factor TU | 405 | С | Elongation Factor TS | 196 |
| | 9 | 1AR1 | 2.70 | В | Cytochrome Oxidase | 298 | А | Cytochrome Oxidase | 558 |
| | 10 | 1AZS | 2.30 | А | VC1 | 220 | В | IIC2 | 212 |
| | 11 | 1AZS | 2.30 | С | GS-ALPHA | 402 | В | IIC2 | 212 |
| | 12 | 1AZZ | 2.30 | С | Ecotin | 142 | А | Collogenase | 226 |
| | 13 | 1B6C | 2.60 | D | TGF-B Superfamily Receptor Type 1 | 342 | С | FK506-Binding Protein | 107 |
| | 14 | 1B7T | 2.50 | Y | Myosin Regulatory Light Chain | 156 | А | Myosin Heavy Chain | 835 |
| | 15 | 1B7T | 2.50 | Z | Myosin Essential Light Chain | 156 | А | Myosin Heavy Chain | 835 |
| | 16 | 1BCC | 3.16 | E | Ubiquinol Cytochrome C Oxidoreductase | 196 | D | Ubiquinol Cytochrome C Oxidoreductase | 241 |
| | 17 | 1BCC | 3.16 | G | Ubiquinol Cytochrome C Oxidoreductase | 81 | С | Ubiquinol Cytochrome C Oxidoreductase | 380 |
| | 18 | 1BGY | 3.00 | R | Cytochrome BC1 Complex | 110 | 0 | Cytochrome BC1 Complex | 379 |
| | 19 | 1BI7 | 3.40 | А | Cyclin-Dependent Kinase 6 | 326 | В | Multiple Tumor Supressor | 156 |
| | 20 | 1BMF | 2.85 | С | Bovine Mitochondrial F1-Atpase | 510 | D | Bovine Mitochondrial F1-Atpase | 482 |
| | 21 | 1BML | 2.90 | А | Plasmin | 250 | С | Streptokinase | 362 |
| | 22 | 1CD1 | 2.67 | А | CD1 | 315 | В | CD1 | 99 |
| _ | 23 | 1CD9 | 2.80 | В | Protein (G-CSF Receptor) | 215 | A | Protein (Granulocyte Colony- Stimulating Factor) | 175 |
| * | 10000 | | | | | | | | |

List of PP complexes from the PIBASE dataset were grouped into 101 heterodimer complexes (**Table B**).

| 24 | 1CN4 | 2.80 | А | Protein (Erythroprotein Receptor) | 228 | С |
|----|------|------|---|--|-----|---|
| 25 | 1DN0 | 2.28 | В | IGM-Kappa Cold Agglutinin (Heavy Chain) | 232 | А |
| 26 | 1DN0 | 2.10 | В | IGM-Kappa Cold Agglutinin (Heavy Chain) | 232 | А |
| 27 | 1DPJ | 1.80 | А | Proteinase A | 329 | В |
| 28 | 1EOC | 2.25 | В | Protocatechuate 3,4-Dioxygenase Beta Chain | 241 | A |
| 29 | 1EP1 | 2.20 | А | Dihydroorotate Dehydrogenase B (Pyrd Subunit) | 311 | В |
| 30 | 1EYS | 2.20 | Μ | Photosynthetic Reaction Center | 324 | С |
| 31 | 1EZV | 2.30 | А | Ubiquinol-Cytochrome C Reductase Complex Core Protein I | 430 | В |
| 32 | 1EZV | 2.30 | В | Ubiquinol-Cytochrome C Reductase Complex Core Protein 2 | 352 | А |
| 33 | 1EZV | 2.30 | С | Cytochrome B | 385 | D |
| 34 | 1EZV | 2.30 | D | Cytochrome C1 | 245 | С |
| 35 | 1F60 | 1.67 | А | Elongation Factor EEF1A | 458 | В |
| 36 | 1F8U | 2.90 | А | Acetylcholinesterase | 583 | В |
| 37 | 1FJG | 3.00 | E | 30s Ribosomal Protein S5 | 162 | Н |
| 38 | 1FRT | 4.50 | А | Neonatal FC Receptor | 269 | В |
| 39 | 1FSK | 2.90 | J | Major Pollen Allergen Bet V 1-A | 159 | L |
| 40 | 1FVU | 1.80 | В | Botrocetin Beta Chain | 125 | А |
| 41 | 1FX0 | 3.20 | В | ATP Synthase Beta Chain | 498 | А |
| 42 | 1G73 | 2.00 | С | Inhibitors Of Apoptosis-Like Protein Ilp | 121 | А |
| 43 | 1G9M | 2.20 | С | T-Cell Surface Glycoprotein Cd4 | 185 | G |
| 44 | 1GPW | 2.40 | А | Hisf Protein | 253 | В |
| 45 | 1HR6 | 2.50 | D | Mitochondrial Processing Peptidase Beta Subunit | 443 | С |

| Protein (Erythroprotein) | 166 |
|--|-----|
| IGM-Kappa Cold Agglutinin (Light Chain) | 215 |
| IGM-Kappa Cold Agglutinin (Light Chain) | 215 |
| Proteinase Inhibitor IA3 Peptide | 33 |
| Protocatechuate 3,4-Dioxygenase Alpha Chain | 209 |
| Dihydroorotate Dehydrogenase B (Pyrk Subunit) | 261 |
| Photosynthetic Reaction Center | 382 |
| Ubiquinol-Cytochrome C Reductase Complex Core Protein 2 | 352 |
| Ubiquinol-Cytochrome C Reductase Complex Core Protein I | 430 |
| Cytochrome C1 | 245 |
| Cytochrome B | 385 |
| Elongation Factor EEF1BA | 94 |
| Fasciculin II | 61 |
| 30s Ribosomal Protein S8 | 138 |
| Beta 2-Microglobulin | 99 |
| Immunoglobulin Kappa Light Chain | 214 |
| Botrocetin Alpha Chain | 133 |
| Atp Synthase Alpha Chain | 507 |
| Second Mitochondria-Derived Activator Of Caspases | 162 |
| Envelope Glycoprotein GP120 | 321 |
| Amidotransferase HISH | 201 |
| Mitochondrial Processing Peptidase Alpha Subunit | 475 |
| | |
| 46 47 | 1HXM 1HXM | 3.12 3.12 | A B | Gamma-Delta T-Cell Receptor Gamma-Delta T-Cell Receptor | 229 242 | B A |
|----------|--------------|--------------|--------|--|------------|--------|
| 48 | 1HYR | 2.70 | А | NKG2-D Type li Integral Membrane Protein | 137 | С |
| 49 | 1I1R | 2.40 | А | Interleukin-6 Receptor Beta Chain Fusion Protein Consisting Of Mhc E- | 301 | В |
| 50 | 1I3R | 2.40 | D | Beta-K Precursor, Glycine Rich Linker, And Hemoglobin Beta-2 Chain | 228 | С |
| 51 | 1185 | 3.20 | D | Cytotoxic T-Lymphocyte-Associated Protein 4 | 110 | В |
| 52 | 1JEB | 2.10 | А | Hemoglobin Zeta Chain | 142 | D |
| 53 | 1JFF | Х | А | Tubulin Alpha Chain | 451 | В |
| 54 | 1JWH | 3.10 | А | Casein Kinase II, Alpha Chain | 337 | D |
| 55 | 1JWI | 2.00 | В | Platelet Aggregation Inducer | 125 | А |
| 56 | 1KFY | 3.60 | В | Fumarate Reductase Iron-Sulfur Protein | 243 | С |
| 57 | 1KFY | 3.60 | С | Fumarate Reductase 15 Kda Hydrophobic Protein | 130 | В |
| 58 | 1KFY | 3.60 | D | Fumarate Reductase 13 Kda Hydrophobic Protein | 119 | В |
| 59 | 1KSG | 2.30 | А | Arf-Like Protein 2 | 186 | В |
| 60 | 1I W6 | 1 50 | F | Subtilisin BPN | 281 | Т |
| 61 | 1M56 | 2 30 | Δ | Cytochrome C Oxidase | 566 | B |
| 01 | 11000 | 2.00 | / \ | | 000 | |
| 62 | 1MF8 | 3.10 | В | Calcineurin B Subunit Isoform 1 | 170 | A |
| 63 | 1MG2 | 2.25 | 0 | Amicyanin | 105 | Ν |
| 64 | 1NCC | 2.50 | Ν | Influenza A Subtype N9 Neuraminidase | 389 | Н |
| 65 | 1NF3 | 2.10 | В | G25K GTP-Binding Protein, Placental Isoform | 195 | D |
| - | | | | | | |

| Gamma-Delta T-Cell Receptor Gamma-Delta T-Cell Receptor | 242 229 |
|---|--------------------------|
| MHC Class I Chain-Related Protein A | 275 |
| Viral IL-6 | 167 |
| H-2 Class II Histocompatibility Antigen, E-K Alpha Chain | 192 |
| T Lymphocyte Activation Antigen CD86 | 126 |
| Hemoglobin Beta-Single Chain Tubulin Beta Chain Casein Kinase II Beta Chain Bitiscetin | 146 445 215 131 |
| Fumarate Reductase 15 Kda | 130 |
| Fumarate Reductase Iron-Sulfur Protein | 243 |
| Fumarate Reductase Iron-Sulfur Protein | 243 |
| Retinal Rod Rhodopsin-Sensitive Cgmp 3',5'-Cyclic Phosphodiesterase Delta-Subunit | 152 |
| Ubtilisin-Chymotrypsin Inhibitor-2A Cytochrome C Oxidase | 63 264 |
| Calmodulin-Dependent Calcineurin A | 373 |
| Methylamine Dehydrogenase, Light Chain | 131 |
| IGG2A-Kappa NC41 FAB (Heavy Chain) | 221 |
| PAR-6B | 128 |

| 66 | 1NFD | 2.80 | А | N15 Alpha-Beta T-Cell Receptor | 203 | D |
|-------|------|------|---|--|-----|---|
| 67 | 105D | 2.05 | L | Coagulation Factor VII | 152 | Т |
| 68 | 10GA | 1.40 | D | T-Cell Receptor Alpha Chain V Region | 215 | Е |
| 69 | 10GA | 1.40 | E | T-Cell Receptor Beta Chain C Region | 252 | D |
| 70 | 1P0S | 2.80 | E | Ecotin Precursor | 142 | Н |
| 71 | 1P4B | 2.35 | L | Antibody Variable Light Chain | 135 | Н |
| 72 | 1PKQ | 3.00 | А | (8-18C5) Chimeric Fab, Light Chain | 241 | В |
| 73 | 1PYT | 2.35 | В | Procarboxypeptidase A | 309 | А |
| 74 | 1Q90 | 3.10 | В | Cytochrome B6 | 215 | D |
| 75 | 1QA9 | 3.20 | А | Human Cd2 Protein | 102 | D |
| 76 | 1R4P | 1.77 | Е | Shiga-Like Toxin Type II B Subunit | 70 | А |
| 77 | 1R8Q | 1.86 | E | Arno | 203 | А |
| 78 | 1SB2 | 1.90 | А | Rhodocetin Alpha Subunit | 133 | В |
| 79 | 1SPG | 1.95 | В | Hemoglobin | 147 | А |
| 80 | 1SQB | 2.69 | В | Ubiquinol-Cytochrome C Reductase Complex Core Protein 2, Mitochondrial | 453 | Ι |
| 81 | 1SQP | 2.70 | A | Ubiquinol-Cytochrome-C Reductase Complex Core Protein I, Mitochondrial Precursor | 480 | I |
| 82 | 1SVX | 2.24 | В | Maltose-Binding Periplasmic Protein | 395 | А |
| 83 | 1UKM | 1.90 | А | EMS16 A Chain | 134 | В |
| Conti | nued | | | | | |

| N15 Alpha-Beta T-Cell Receptor Tissue factor | 239 218 |
|---|---|
| T-Cell Receptor Beta Chain C Region | 252 |
| T-Cell Receptor Alpha Chain V Region Coagulation Factor X Precursor Antibody Variable Heavy Chain (8-18C5) Chimeric Fab, Heavy Chain Procarboxypeptidase A Cytochrome B6-F Complex Subunit 4 Human CD58 Protein Shiga-Like Toxin Type II A Subunit ADP-Ribosylation Factor 1 Rhodocetin Beta Subunit Hemoglobin | 215 254 124 252 94 159 95 297 181 129 144 |
| Ubiquinol-Cytochrome C Reductase 8 Kda Protein | 78 |
| Ubiquinol-Cytochrome C Reductase Iron-Sulfur Subunit, Mitochondrial Precursor (EC 1.10.2.2) (Rieske Iron- Sulfur Protein) (RISP) [Contains: Ubiquinol-Cytochrome C Reductase 8 Kda Protein (Complex III Subunit IX)] | 78 |
| Ankyrin Repeat Protein Off7 | 169 |
| EMS16 B Chain | 128 |

| 84 | 1UKV | 1.50 | G | Secretory Pathway GDP Dissociation Inhibitor | 453 | Y | GTP-Binding Protein YPT1 | 206 |
|-----|------|------|---|---|-----|---|---|-----|
| 85 | 1UVQ | 1.80 | А | HLA Class II Histocompatibility Antigen | 197 | В | Hla Class II Histocompatibility Antigen | 198 |
| 86 | 1VF5 | 3.00 | А | Cytochrome B6 | 215 | В | Subunit IV | 160 |
| 87 | 1X9F | 2.60 | J | Globin II, Extracellular | 145 | K | Globin III, Extracellular | 153 |
| 88 | 1XDK | 2.90 | В | Retinoic Acid Receptor, Beta | 303 | А | Retinoic Acid Receptor RXR-Alpha | 238 |
| 89 | 1XXD | 2.91 | В | Coagulation Factor XI | 238 | С | Ecotin | 142 |
| 90 | 1YTZ | 3.00 | С | Troponin C | 162 | I | Troponin I | 182 |
| 91 | 1YVB | 2.70 | А | Falcipain 2 | 241 | I | Cystatin | 111 |
| 92 | 1Z7M | 2.90 | E | ATP Phosphoribosyltransferase | 208 | А | ATP Phosphoribosyltransferase Regulatory Subunit | 344 |
| 93 | 1Z7X | 1.95 | Х | Ribonuclease I | 129 | W | Ribonuclease Inhibitor | 461 |
| 94 | 2A9K | 1.73 | А | Ras-Related Protein Ral-A | 187 | В | Mono-ADP-Ribosyltransferase C3 | 223 |
| 95 | 2A9K | 1.73 | В | Mono-ADP-Ribosyltransferase C3 | 223 | А | Ras-Related Protein Ral-A | 187 |
| 96 | 2BEX | 1.99 | С | Nonsecretory Ribonuclease | 135 | А | Ribonuclease Inhibitor | 460 |
| 97 | 2G2U | 1.60 | А | Beta-lactamase SHV-1 | 265 | В | Beta-Lactamase Inhibitory Protein | 165 |
| 98 | 2GJ7 | 5.00 | E | Glycoprotein E | 401 | В | IG Gamma-1 Chain C Region | 227 |
| 99 | 2IG2 | 3.00 | н | IGG1-Lambda Kol FAB (Heavy Chain) | 455 | L | IGG1-Lambda KOL FAB (Light Chain) | 216 |
| 100 | 2J12 | 1.50 | В | Coxsackievirus And Adenovirus Receptor | 128 | А | Fiber Protein | 194 |
| 101 | 20NL | 4.00 | С | MAP Kinase-Activated Protein Kinase 2 | 406 | А | Mitogen-Activated Protein Kinase 14 | 366 |

8.2 Supplementary Information for Chapter 4

Table 4.1 List of the transcription factors targeted to STIM and ORAI genes with Pubmed ID (PMID), technique of experiments and cell types obtained from the CheA database.

| | | PMID | | |
|-----------|---------------|----------------------|----------------------|---|
| Gene | TF | (Pubmed) | Technique | Cell Туре |
| STIM1 | AR | 20517297 | CHIP-SEQ | VCAP |
| | E2F4 | 21247883 | CHIP-SEQ | LYMPHOBLASTOID |
| | EGR1 | 20690147 | CHIP-SEQ | ERYTHROLEUKEMIA |
| | ELK3 | 25401928 | CHIP-SEQ | HUVEC |
| | FLI1 | 21571218 | CHIP-SEQ | MEGAKARYOCYTES |
| | FOXA2 | 19822575 | CHIP-SEQ | HepG2 |
| | FOXP1 | 21924763 | CHIP-SEQ | HESC |
| | GATA1 | 21571218 | CHIP-SEQ | MEGAKARYOCYTES |
| | GATA1 | 19941826 | CHIP-SEQ | K562 |
| | GATA2 | 19941826 | CHIP-SEQ | K562 |
| | GATA2 | 21666600 | CHIP-SEQ | HMVEC |
| | HNF4A | 19822575 | CHIP-SEQ | HepG2 |
| | MITF | 21258399 | CHIP-SEQ | MELANOMA |
| | MYC | 19915707 | CHIP-SEQ | AK7 |
| | NCOR1 | 26117541 | CHIP-SEQ | K562 |
| | PHF8 | 20622854 | CHIP-SEQ | HELA |
| | RUNX1 | 17652178 | CHIP-SEQ | JURKAT |
| | SCL | 21571218 | CHIP-SEQ | MEGAKARYOCYTES |
| | SOX2 | 21211035 | CHIP-SEQ | LN229_GBM |
| | SPI1 | 20517297 | CHIP-SEQ | HL60 |
| | TP63 | 22573176 | CHIP-SEQ | HFKS |
| | TRIM28 | 17542650 | CHIP-SEQ | NTERA2 |
| | TTF2 | 22483619 | CHIP-SEQ | HELA |
| STIM2 | AR | 22383394 | CHIP-SEQ | PROSTATE_CANCER |
| | AR | 19668381 | CHIP-SEQ | PC3 |
| | AR CUX1 | 25329375 19635798 | CHIP-SEQ CHIP-SEQ | VCAP MULTIPLE HUMAN CANCER CELL TYPES |
| | E2F1 | 21310950 | CHIP-SEQ | MCF7 |
| | E2F4 | 17652178 | CHIP-SEQ | JURKAT |
| | EGR1 | 20690147 | CHIP-SEQ | ERYTHROLEUKEMIA |
| | ELK3 | 25401928 | CHIP-SEQ | HUVEC |
| | FOXP1 | 21924763 | CHIP-SEQ | HESC |
| | GABP | 19822575 | CHIP-SEQ | HepG2 |
| | GATA1 | 21571218 | CHIP-SEQ | MEGAKARYOCYTES |
| | GATA2 | 19941826 | CHIP-SEQ | K562 |
| | HOXB7 | 26014856 | CHIP-SEQ | BT474 |
| | MYCN | 21190229 | CHIP-SEQ | SHEP-21N |
| | PAX3- EKHP | 20663000 | | RHABDOMYOSARCOMA |
| Continued | | 20003003 | | |

| | POU3F2 | 20337985 | CHIP-SEQ | 501MEL |
|----------|--------|----------|----------|--------------------|
| | PPARD | 21283829 | CHIP-SEQ | MYOFIBROBLAST |
| | RUNX1 | 17652178 | CHIP-SEQ | JURKAT |
| | RUNX2 | 22187159 | CHIP-SEQ | PCA |
| | SCL | 21571218 | CHIP-SEQ | MEGAKARYOCYTES |
| | SMAD4 | 21799915 | CHIP-SEQ | A2780 |
| | SOX11 | 23321250 | CHIP-SEQ | Z138-A519-JVM2 |
| | SOX2 | 21211035 | CHIP-SEQ | LN229_GBM |
| | TOP2B | 26459242 | CHIP-SEQ | MCF7 |
| | TP63 | 23658742 | CHIP-SEQ | EP156T |
| | TTF2 | 22483619 | CHIP-SEQ | HELA |
| | VDR | 23849224 | CHIP-SEQ | CD4+ |
| | WT1 | 25993318 | CHIP-SEQ | PODOCYTE |
| | ZNF217 | 24962896 | CHIP-SEQ | MCF7 |
| ORAI1 | BCL6 | 25482012 | CHIP-SEQ | CML/JURL-MK1 |
| | E2F4 | 21247883 | CHIP-SEQ | LYMPHOBLASTOID |
| | ELK3 | 25401928 | CHIP-SEQ | HUVEC |
| | FLI1 | 21571218 | CHIP-SEQ | MEGAKARYOCYTES |
| | FOXA2 | 19822575 | CHIP-SEQ | HepG2 |
| | FOXP1 | 21924763 | CHIP-SEQ | HESC |
| | FOXP2 | 23625967 | CHIP-SEQ | PFSK-1 AND SK-N-MC |
| | GATA6 | 25053715 | CHIP-SEQ | YYC3 |
| | KLF5 | 25053715 | CHIP-SEQ | YYC3 |
| | MITF | 21258399 | CHIP-SEQ | MELANOMA |
| | MYC | 22102868 | CHIP-SEQ | BL |
| | PPARD | 21283829 | CHIP-SEQ | MYOFIBROBLAST |
| | RUNX1 | 21571218 | CHIP-SEQ | MEGAKARYOCYTES |
| | SMAD4 | 21741376 | CHIP-SEQ | HESC |
| | TFEB | 21752829 | CHIP-SEQ | HELA |
| ORAI2 | EBNA2 | 21746931 | CHIP-SEQ | IB4-LCL |
| | ELK3 | 25401928 | CHIP-SEQ | HUVEC |
| | GATA1 | 19941826 | CHIP-SEQ | K562 |
| | GATA2 | 19941826 | CHIP-SEQ | K562 |
| | MITF | 21258399 | CHIP-SEQ | MELANOMA |
| | MYC | 22102868 | CHIP-SEQ | BL |
| | PPARD | 21283829 | CHIP-SEQ | MYOFIBROBLAST |
| | RBPJ | 21746931 | CHIP-SEQ | IB4-LCL |
| | RUNX1 | 21571218 | CHIP-SEQ | MEGAKARYOCYTES |
| | TP63 | 22573176 | CHIP-SEQ | HFKS |
| | VDR | 24763502 | CHIP-SEQ | THP-1 |
| ORAI3 | BCL6 | 25482012 | CHIP-SEQ | CML/JURL-MK1 |
| | ELF1 | 20517297 | CHIP-SEQ | JURKAT |
| | ELK3 | 25401928 | CHIP-SEQ | HUVEC |
| | FOXA1 | 25552417 | CHIP-SEQ | VCAP |
| (,• 1 | GATA2 | 19941826 | CHIP-SEQ | K562 |

| GATA2 | 21666600 | CHIP-SEQ | HMVEC | |
|--------|----------|----------|---------------|--|
| GATA3 | 24758297 | CHIP-SEQ | MCF7 | |
| GATA4 | 25053715 | CHIP-SEQ | YYC3 | |
| HNF4A | 19822575 | CHIP-SEQ | HepG2 | |
| MITF | 21258399 | CHIP-SEQ | MELANOMA | |
| NCOR1 | 26117541 | CHIP-SEQ | K562 | |
| PPARD | 21283829 | CHIP-SEQ | MYOFIBROBLAST | |
| SPI1 | 23547873 | CHIP-SEQ | NB4 | |
| SPI1 | 23127762 | CHIP-SEQ | K562 | |
| TFAP2C | 20629094 | CHIP-SEQ | MCF7 | |

TF, transcription factor.

| Promoter Name | Promoter Sequence |
|------------------|---|
| STIM1_1 | TTGGGGGCTGGGAGCTCGCCCCCGGGCCGAGCCGGGTCAGGCTGTTGTCGCCTCAGGCAG |
| | CTCCTGGGAGGCTAACGTCGTGTCCTGGGCCTCTGTTTAGACAGCTCTAGAACTGAGGCG |
| | AGTGGAGCAGCACCAAGGCCCGGAGATCGGGGGCAGGGCAGCTGCTGTCGCCGCCGCCGCA |
| | GGCCTGAGTTACCTGAGTAACTGCGGGTCAGGGACCCGCCGACGGCCCGCGGTTGGCGC |
| | TGGAGACTCTCGGTGGGGAAAGGGAAGCTGGGACTTGATCCTTTGCGCGGGATCCTGGCA |
| | AAGACTAGCGCGGGGCCGGGGGGTCCGGGAGAGCCCGCTAGGGGCGGGGGATTCCGGGGAGCC |
| | GTCTTCACCGGTTATTCCGGGATCCAGCTGGGCGCTGGGGCTGGCCCGGGCTTCGCTGGG |
| | GACCGGGCGGGGGGGGGGGGGGGGGGGGGGGGGGCGCGCCCGCCCGGCCCGGGCCCGCCCC |
| | GCGCCGCCCGCCTGGAAGCCGCTGTCCTGGGCCTGGCCGGTGTGCGTCCGCCTGCT |
| | GGACCTGGGCACCGCCAGCCGCCTGGGCACGGGACTGGGCGGGGGGCGCTGACCTCGGCCT |
| STIM1_2 | CTCATTTCTTTCATTGCTTACAAGTAGATAGCATTCCAGTTCATAGGTTTCTTTGAGAAA |
| | TAGACTGTAGAAAAGACGACAATGTTTATTCTCATTAGTCAGACGAACTGCAGCACAAGT |
| | GTTTCAAAGAGGAACCCCCAAGAACTCCCAGGCTTGTAGGAAGCATCTGATTTTACATAAG |
| | TTTGCACAGTAGGAAACTGAGGTCCTCGGGGAGGAACAAAGCAACTTATGATCAGACAAA |
| | TGAGTCACTAGTAGAGCTGAAATGAGAAGCCAGATCTCCTAACCACTCCCACACCCCATC |
| | ACAGTGCCTCACTGCGTCTCTTACTGGTGGGCTTGATTGCTTTCCAAGGCCAGAGAAGGA |
| | AGTAGCTGTTCCTGTATGCTCAGACAGGAGTATAAATCACACTGTGATGTCAGAAGCTTC |
| | TTTTCTAGCTGGAGAAATAAAGCTATGCATCAGAAAAGAGCACCAATCTTATAGAGTAGT |
| | ATAGAATTAAGTGCTTACTTGTGGGACCTAGACCAGTGGTTTTCAACCTTTCATATCTTT |
| | TACCCCCCTCCTCTTATATCCACTGTGTTCCAGCAGAAGAGAAACAAAGTTCCATTTTC |
| STIM2_1 | AAAACGCTAAGCTATGCTAACCGCGTCTAAACAGCCAGCC |
| | TTGAAATTACGGGGACGCTGCTGTTTTCCGATTGCAAGATTTTCCTGATCCATGCAATTA |
| | CTTTCGCTGCCCTCTACGAGGCTGAAACTCGCCCTCAGGATGTGGGACGTCTGGACTCTT |
| | CTCTCCGTCCCCTTGTAGCCCCCCACTCCCCCTCGCGGTGGTACCGTGAATGAGGGAGAG |
| | GTACACGTCCCCCTTCTTCCCCGCCTCCTATCTTCGCGGCTCGCTAAAGCGTTATCAGCC |
| | GCCCCACGGTACTACCGTCCGTCTAGGAACGCCTCCGGGGCGGGGCTGGGATGCCGCGCA |
| | CGCGCAGTACAGCAGCGCCGCGCCTGCGCCGTGGAGAGCCTGAGGGAGG |
| | TATGCGAGCGAATGTGCGAGGGGGGGGGGGGGGGGGCGTCCCGGCGGGGGGGG |
| | CGCGGGCCGGAGGGGGGGGGGGGGGGGGCGCGGCGGGGGG |
| | GTGTTTGGCGGCGCCAGAGCAGCGGATCCCGGTCTCGCCGCAGCAGCAGCGCGGGTGTCG |
| STIM2_2 | TACAGATTTATTATAAATGTGTATACTTGAGAATAAATGAGTAGAAGAACAGTGGAGGTA |
| | AGTCAAATATAGTGGAATTAGATGTGTAGTTAAATTTTATTTTTAACTTGATTTAAATAA |
| | TTAGGAATTTTTGAAAAGCTTTTTGCGGAAGAGTATTTCCCTGCTTTCCCTGTCATTTGA |
| | ACCCAGGATAACCAAAATAGCTGTATAGTAAGTTGCCTGATATTTGTATTAACCAAACTT |
| | AAGGCTAATGAAAAATGCTATGATTTCTGATTGAAATATGTATTTAATCGCTTGACCCAG |
| | TATTCATTATTTTGTAAAAAAAAAAAAAACTGGAAATTTTTGTGAGGAATTTTTATTTTT |
| | ATTGTTCTATAAGGCTTGAAAAGGCACAGGAAGAAAACAGAAATGTTGCTGTAGAAAAGC |
| | AAAATTTAGAGCGCAAAATGATGGATGAAATCAATTATGCAAAGGAGGAGGCTTGTCGGC |
| | TGAGAGAGCTAAGGGAGGGAGCTGAATGTGAATTGAGTAGACGTCAGTATGCAGAACAGG |
| | AATTGGAACAGGTATTTACATTAAAAAAAAAAAACCACTTGTAAAGATGTTAACATTGCCAC |
| Continued | 1 |

Table 4.2 List of the promoter sequences of STIM and ORAI genes obtained from the EPDnew database.

| ORAI1 | GACAGGTTTTGTCCATCTCGTTCACCACCCTACCCAGCCTCAGCACCTAGACCAGTGTTG |
|-------|--|
| | GCACCCAGTGGGCGCCAAATAAACACTGCTTGAACTCCAGACGTCAGCCGCTCTTTTTCC |
| | TACAGACCTTGAGCCACCTTGTTCCAAAGGGGATATGGGCCTCAGGAGGCGCCCAGAGGT |
| | GACCTCAGGCGGCCCGACCCAGGAGTCCAAGCTCCAGGAGCAGGGCCACGGGAGCAGCTG |
| | CGGAGAGGGGGGGGCGCCAGGAGCCGGAGCGGGCAGCCGGGCGCTTCCAGGAAAAGTGGCG |
| | GGCGGCGCGCGCCAGGGACCGTGGGCGGTGCCGTCGGAGCGGGCGG |
| | CACAACAACGCCCACTTCTTGGTGGGCGGGGGCACAGGTGGGCGGGGGGGAGCATGCAAAACAG |
| | CCCAGGGCGGCGGCCAATCGCGGGGCGCGCGGGGGTCCAGGCCCCGGGGATCCGAGGCG |
| | CCGCCCGCGCGCAGTCTCTGGTCACTGCCGCCCGGGGGCTTTTGCCAGCGGCGCGCGGG |
| | CCTGCGTGCTGGGGCAGCGGGCACTTCTTCGACCTCGTCCTCCTCGTCCTGTGCGGCCGG |
| ORAI2 | ACTCCCCGCCCTCTACGCAGCCAGCGTCCAATGCTGGGCACCCCCGAGGCTCACCCTG |
| | CCAAGCCTGGGGCTCCCCTTTTGCGCCCGGACCCAGGGGCAGGGAAAGCCCAGCTCGTGG |
| | TCTGTGGGTAGCCGGACCCCCGATGGGGGGGGGGGGGCCTCGCCTTGACTCCCAGAGCT |
| | GGGGCCGGGGACAGGAGCTGGGGCAGGAGGGATGCGCGCGGGTCGGGGTCTTCCCACCTC |
| | CCCTGCTCCTCCCCCGCGATCCCGGGGTGGTTCCAGGTGAGGCGGGGACCCCCACC |
| | CCCCCACTCTCCGAGGAGGCGCCGCCAGCCCGCCCCCCCC |
| | CCGCGGCTCTCCCGCGGGTGGGTCACGTGTTGGCGGCGCCTGGTTGCCTTGGCAGCGGCT |
| | GCGGCGGCCGCGGGGGGGGGGGGGGGGGGGGGGGGGGGG |
| | GAGGGAGGGGGCCGCGCTCGGCGCCCGGGCCACTGGGCCACAGGCCACGCGGCCA |
| | CGCAGTCCGAGCGGGAGCCGAGCCGGGGGGGGGGGGGGG |
| ORAI3 | GTAACAGGGAGGTGCGCGGGGGGGGGGGGGGGGGGGGGG |
| | GCCTGGGGATAGCGAGAGGCTTGAGAATGGGGCCGCTTGGGGGAGGGA |
| | GCGAGGGGCAAGCGGGGGGCCCAGCCGGGCTGGGCCCCTGGGCCCCGGGTCTGTACAATA |
| | CGGTTTGCTATAAAACTCAAAATCTTCCAGCCGGGGCTGCGGAGTTCGTGTGTGT |
| | CGGGGTCCCTACCTACAGATGAGTGGGCTCACCTCTCCTGGACTCATTTTGGGAGGGA |
| | TGGAAGTGTGGACACCTGGGGTGTCCAGCTGTACCTTGGAGGGGGCTGGGGTTGGCGTGC |
| | ACCTCGGTGGGGTCCGGGCGCTTGGATAACGTTCTTGGTGGGTAGGGGTCGCGGGGAATC |
| | TCTGCGGGCCCGGGACTGCGGGGACTTGGTCCCCGGCTCCACCCCATCATGTGGCTAGCC |
| | CCGGCTCCGCCTCTGTCCCAGTTCCTGTTTTGGCCTCCGCTGTCCCGCTCCGGCTCCTGG |
| | GGCTCCCCGCAGACGCTGCTTTTCTTGCTCCACTGGGGGTGCCTCTTCCTGGGCGCCCCGC |

| Promoter name | Potential Interactions | Type of interaction | Bridge Proteins | Interactions | node1 | node2 | Score | node1 | node2 | Score | node1 | node2 | Score |
|------------------|---------------------------|------------------------|--------------------|----------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| STIM2_1 | E2F1:E2F4 | Direct | | | E2F1 | E2F4 | 0.966 | | | | | | |
| | E2F1:E2F4 | Bridge protein | MYC | E2F1>MYC>E2F4 | E2F1 | MYC | 0.982 | MYC | E2F4 | 0.982 | | | |
| STIM2_1 | E2F1:EGR1 | Bridge protein | SPI1 | E2F1>SPI1>EGR1 | E2F1 | SPI1 | 0.905 | SPI1 | EGR1 | 0.501 | | | |
| | E2F1:EGR1 | Bridge protein | E2F4 | E2F1>E2F4>EGR1 | E2F1 | E2F4 | 0.966 | E2F4 | EGR1 | 0.846 | | | |
| | E2F1:EGR1 | Bridge protein | MYC | E2F1>MYC>EGR1 | E2F1 | MYC | 0.982 | MYC | EGR1 | 0.928 | | | |
| STIM2_1 | E2F1:GATA1 | Bridge protein | ELF1 | E2F1>ELF1>GATA1 | E2F1 | ELF1 | 0.902 | ELF1 | GATA1 | 0.893 | | | |
| | E2F1:GATA1 | Bridge protein | SPI1 | E2F1>SPI1>GATA1 | E2F1 | SPI1 | 0.905 | SPI1 | GATA1 | 0.995 | | | |
| | E2F1:GATA1 | Bridge protein | MYC | E2F1>MYC>GATA1 | E2F1 | MYC | 0.982 | MYC | GATA1 | 0.961 | | | |
| STIM2_1 | E2F4:GATA1 | Direct | | | E2F4 | GATA1 | 0.727 | | | | | | |
| | E2F4:GATA1 | Bridge protein | MYC | E2F4>MYC>GATA1 | E2F4 | MYC | 0.96 | MYC | GATA1 | 0.961 | | | |
| | E2F4:GATA1 | Bridge protein | EGR1 | E2F4>EGR1>GATA1 | E2F4 | EGR1 | 0.846 | EGR1 | GATA1 | 0.584 | | | |
| ORAI1 | PPARD:E2F4 | Bridge protein | EGR1 | PPARD>EGR1>E2F4 | PPARD | EGR1 | 0.691 | EGR1 | E2F4 | 0.846 | | | |
| ORAI2 | NFKB1:RUNX1 | Bridge protein | EGR1 | NFKB1>EGR1>RUNX1 | NFKB1 | EGR1 | 0.671 | EGR1 | RUNX1 | 0.676 | | | |
| | NFKB1:RUNX1 | Bridge protein | MYC | NFKB1>MYC>RUNX1 | NFKB1 | MYC | 0.993 | MYC | RUNX1 | 0.61 | | | |
| | NFKB1:RUNX1 | Bridge protein | SPI1 | NFKB1>SPI1>RUNX1 | NFKB1 | SPI1 | 0.781 | SPI1 | RUNX1 | 0.991 | | | |
| | PPARD:MYC | Bridge protein | EGR1 | PPARD>EGR1>MYC | PPARD | EGR1 | 0.691 | EGR1 | MYC | 0.928 | | | |
| | PPARD:MYC | Bridge protein | HNF4A | PPARD>HNF4A>MYC | PPARD | HNF4A | 0.916 | HNF4A | MYC | 0.942 | | | |
| ORAI3 | MITF:ELF1 | Bridge proteins | SPI1, E2F1 | MITF>SPI1>E2F1>ELF1 | MITF | SPI1 | 0.951 | SPI1 | E2F1 | 0.905 | E2F1 | ELF1 | 0.902 |
| | MITF:ELF1 | Bridge proteins | SPI1, GATA1 | MITF>SPI1>GATA1>ELF1 | MITF | SPI1 | 0.951 | SPI1 | GATA1 | 0.995 | GATA1 | ELF1 | 0.893 |
| | MITF:SPI1 | Direct | | | MITF | SPI1 | 0.951 | | | | | | |
| | SPI1:ELF1 | Bridge protein | GATA1 | SPI1>GATA1>ELF11 | SPI1 | GATA1 | 0.995 | GATA1 | ELF1 | 0.893 | | | |
| | SPI1:ELF1 | Bridge protein | E2F1 | SPI1>E2F1>ELF1 | SPI1 | E2F1 | 0.905 | E2F1 | ELF1 | 0.902 | | | |

Table 4.3 List of the predicted interactions between transcription factors, name of bridge proteins and score obtained from the STRING database.

| Node 1 Node 2 | Node 1 | Node 2 | Weight | |
|---------------|--------|--------|-----------|--------|
| Normal-tum | or | Τι | umor-norm | al |
| EP300 CREBBP | 0.2822 | FOXP2 | 0.1546 | |
| NCOR1 EP300 | 0.2343 | MITF | RUNX2 | 0.1528 |
| SP1 EP300 | 0.2093 | FOXA1 | GATA3 | 0.1112 |
| NCOR1 CREBBP | 0.1903 | FLI1 | ELK3 | 0.1095 |
| MITF ELK3 | 0.1636 | SPI1 | FLI1 | 0.1049 |
| NCOR1 AR | 0.1624 | AR | EP300 | 0.0869 |
| SMAD4 SP1 | 0.1448 | SPI1 | TFEB | 0.0749 |
| MITF FLI1 | 0.1331 | FOXP2 | FLI1 | 0.0737 |
| SP1 CREBBP | 0.1330 | SP1 | EP300 | 0.0673 |
| FOP2B SP1 | 0.1270 | FOXP2 | RUNX2 | 0.0625 |
| ELK3 | 0.1135 | FOXP2 | MITF | 0.0623 |
| OXA1 GATA3 | 0.0985 | TOP2B | SP1 | 0.0623 |
| VCOR2 TRIM28 | 0.0976 | STIM2 | FLI1 | 0.0604 |
| MITF TRIM28 | 0.0963 | ZNF217 | TFEB | 0.0584 |
| NCOR1 SP1 | 0.0929 | MITF | ELK3 | 0.0579 |
| AR EP300 | 0.0915 | SP1 | AR | 0.0552 |
| TTF2 EP300 | 0.0881 | FOXA1 | FLI1 | 0.0552 |
| TRIM28 ELK3 | 0.0829 | FLI1 | TFEB | 0.0536 |
| FLI1 HOXB7 | 0.0820 | TOP2B | EP300 | 0.0532 |
| KDM1A HOXB7 | 0.0807 | NCOR1 | SP1 | 0.0487 |
| SMAD4 TOP2B | 0.0805 | AR | CREBBP | 0.0484 |
| KDM1A HDAC2 | 0.0803 | ELF1 | TOP2B | 0.0474 |
| KDM1A FLI1 | 0.0743 | TP63 | KLF5 | 0.0470 |
| FTF2 SP1 | 0.0698 | EP300 | ELK3 | 0.0440 |
| RBPJ TRIM28 | 0.0688 | RUNX2 | ELK3 | 0.0431 |
| KLF5 TFAP2C | 0.0686 | E2F1 | RUNX1 | 0.0424 |
| HOXB7 TFAP2C | 0.0670 | ELF1 | SP1 | 0.0410 |
| HDAC3 KDM1A | 0.0647 | FOXP2 | EGR1 | 0.0387 |
| KDM1A MITF | 0.0624 | TOP2B | AR | 0.0378 |
| TOP2B EP300 | 0.0603 | FLI1 | GATA3 | 0.0378 |
| FOXP1 TP63 | 0.0587 | TP63 | EGR1 | 0.0369 |
| FOXP2 HDAC1 | 0.0549 | EGR1 | ELK3 | 0.0363 |
| AR CREBBP | 0.0548 | EGR1 | FLI1 | 0.0360 |
| MITF VDR | 0.0544 | NCOR1 | EP300 | 0.0349 |
| TRIM28 FLI1 | 0.0539 | SP1 | CREBBP | 0.0324 |
| MITF RBPJ | 0.0539 | STIM2 | ELK3 | 0.0313 |
| FOXP2 MITF | 0.0536 | FOXA1 | E2F4 | 0.0310 |
| HOXB7 VDR | 0.0525 | NFKB1 | ELK3 | 0.0307 |
| MITF HOXB7 | 0.0522 | NCOR1 | AR | 0.0306 |
| FLI1 VDR | 0.0516 | ELF1 | ELK3 | 0.0300 |
| NCOR1 PHF8 | 0.0496 | FOXA1 | SPI1 | 0.0295 |

Table 4.4 List of differential interactions in (normal-tumor) and (tumor-normal) samples.

| SMAD4 | ORAI2 | 0.0491 | MITF | TRIM28 | 0.0294 |
|--------|--------|--------|-------|--------|--------|
| KDM1A | TFAP2C | 0.0486 | SP1 | ELK3 | 0.0283 |
| TP63 | TFAP2C | 0.0478 | MITF | FLI1 | 0.0277 |
| SMAD4 | CUX1 | 0.0457 | PHF8 | CREBBP | 0.0262 |
| TTF2 | TOP2B | 0.0424 | SPI1 | E2F4 | 0.0258 |
| SMAD4 | EP300 | 0.0417 | FOXA1 | TFEB | 0.0257 |
| FOXA1 | PHF8 | 0.0404 | AR | ELK3 | 0.0248 |
| KLF5 | CUX1 | 0.0398 | ELF1 | AR | 0.0247 |
| FOXA1 | AR | 0.0392 | ELF1 | EP300 | 0.0244 |
| ORAI2 | SP1 | 0.0383 | SP1 | CDKN2A | 0.0240 |
| PHF8 | AR | 0.0379 | PHF8 | AR | 0.0239 |
| FLI1 | TFAP2C | 0.0378 | FOXP2 | STIM2 | 0.0235 |
| HOXB7 | GATA2 | 0.0370 | RUNX2 | RUNX1 | 0.0231 |
| KDM1A | ELK3 | 0.0363 | NCOR1 | TOP2B | 0.0225 |
| GATA6 | ELK3 | 0.0360 | NCOR1 | CREBBP | 0.0216 |
| RBPJ | FLI1 | 0.0351 | NFKB1 | AR | 0.0212 |
| FOXP2 | TRIM28 | 0.0331 | FOXP2 | RUNX1 | 0.0211 |
| TTF2 | CREBBP | 0.0328 | ELK3 | RUNX1 | 0.0205 |
| FOXP2 | ELK3 | 0.0306 | ELF1 | NFKB1 | 0.0203 |
| TTF2 | NCOR1 | 0.0302 | FLI1 | GATA6 | 0.0201 |
| KDM1A | VDR | 0.0293 | | | |
| NCOR2 | RBPJ | 0.0286 | | | |
| RBPJ | ELK3 | 0.0282 | | | |
| NCOR1 | TOP2B | 0.0278 | | | |
| TRIM28 | VDR | 0.0277 | | | |
| MITF | TFAP2C | 0.0264 | | | |
| FLI1 | GATA2 | 0.0257 | | | |
| TOP2B | CREBBP | 0.0257 | | | |
| SP1 | AR | 0.0251 | | | |
| MITF | GATA6 | 0.0248 | | | |
| TP63 | CREBBP | 0.0247 | | | |
| KDM1A | TRIM28 | 0.0245 | | | |
| PHF8 | EP300 | 0.0238 | | | |
| VDR | ELK3 | 0.0237 | | | |
| NCOR2 | FLI1 | 0.0236 | | | |
| MITF | NCOR2 | 0.0225 | | | |
| TP63 | HOXB7 | 0.0222 | | | |
| FOXP1 | SMAD4 | 0.0217 | | | |
| HDAC3 | HDAC2 | 0.0211 | | | |
| SP1 | CDKN2A | 0.0207 | | | |
| PHF8 | CREBBP | 0.0201 | | | |
| TOP2B | HOXB7 | 0.0201 | | | |

8.3 Supplementary Information for Chapter 5

| 1Q8T | 1Q8U | 1XD0 | 2YKI | 3JVS | 2J62 | 2W66 | 4DEW |
|------|------|------|------|------|------|------|------|
| 1VSO | 2WBG | 3B3W | 1W3K | 2CBJ | 30ZT | 3GCS | 2X97 |
| 30E5 | 1LOR | 2QBP | 3F3E | 3I3B | 3UO4 | 3ZSX | 2V7A |
| 3L4W | 2QBR | 3PXF | 2YFE | 30V1 | 3HUC | 2WTV | 3F3C |
| 2ZXD | 2VVN | 3F3A | 2D3U | 1H23 | 2ZWZ | 3SU2 | 3BFU |
| 2P4Y | 1YC1 | 2J78 | 2X8Z | 3GBB | 2VW5 | 3GE7 | 1R5Y |
| 4DE2 | 3UDH | 3G2Z | 3EBP | 2GSS | 2XHM | 3MSS | 1U33 |
| 1SLN | 3F17 | 2X00 | 2BRB | 3AO4 | 4DE1 | 4GID | 3D4Z |
| 2ZCR | 3IMC | 10GS | 1GPK | 1PS3 | 3PE2 | 3ACW | 3FCQ |
| 3CYX | 2CET | 3K5V | 2V00 | 1OS0 | 1Z95 | 3BKK | 3MFV |
| 2ZJW | 3DD0 | 3N86 | 3EHY | 3E93 | 1LOQ | 3F80 | 1U1B |
| 3NW9 | 3IVG | 3N7A | 1NVQ | 20BF | 3G0W | 1N2V | 3VH9 |
| 1LBK | 3AG9 | 4G8M | 2G70 | 3SU5 | 2VO5 | 3S8O | 3B68 |
| 3PWW | 2XY9 | 3ZSO | 2FVD | 1KEL | 3UEX | 2ZCQ | 2YGE |
| 3SU3 | 1SQA | 2D10 | 2VOT | 2WEG | 2XNB | 3B3S | 3CJ2 |
| 3FK1 | 4TMN | 4DJR | 1P1Q | 2XDL | 4GQQ | 3GNW | 3DXG |
| 2X0Y | 3MUZ | 1HNN | 3EJR | 2QFT | 2YMD | 2QMJ | 2XYS |
| 3VD4 | 2HB1 | 1HFS | 3L3N | 30WJ | 3L7B | 1W3L | 3COY |
| 3NQ3 | 3NOX | 2ZX6 | 3LKA | 3FV1 | 3MYG | 2PQ9 | 3KV2 |
| 20LE | 1LOL | 1JYQ | 2IWX | 3L4U | 2VL4 | 3KWA | 4DJV |
| 3U9Q | 3UEU | 2WCA | 1E66 | 2XB8 | 3G2N | 1QI0 | |

Table 5.1 Dataset of 167 protein-ligand complexes.

| Program | URL | Year | Limitation |
|----------------|--|-----------|---|
| APROPOS | http://www.csb.yale.edu/poststructure/apropos/apropos.html | 1996 | Need CSB core |
| BetaCavityWeb | http://voronoi.hanyang.ac.kr/betacavityweb/ | 2015 | Available only as webserver |
| CASTp | http://sts.bioe.uic.edu/castp/calculation.php | 2003/2006 | Available only as webserver |
| CAVER Analyst | http://www.caver.cz | 2014 | No output file, not run automatically |
| ConCavity | http://compbio.cs.princeton.edu/concavity/ | 2009 | Not compatible with Debian 16.0-4-amd64 |
| Epock | http://epock.bitbucket.org | 2014 | Need VMD |
| HotPatch | http://hotpatch.mbi.ucla.edu/ | 2007 | Incomplete documentation |
| LIGSITEcs | http://projects.biotec.tu-dresden.de/pocket/ | 1997 | Need BALL algorithm |
| McVol | http://www.bisb.uni-bayreuth.de/index.php?page=data/mcvol/mcvol | 2010 | Does not take PDB as input |
| Metapocket2.0 | http://projects.biotec.tu-dresden.de/metapocket/ | 2011 | Available only as webserver |
| MolSite | http://presto.protein.osaka-u.ac.jp/myPresto4/index.php?lang=en | 2011 | Incomplete documentation |
| PASS | http://www.ccl.net/cca/software/UNIX/pass/overview.shtml | 2000 | Not compatible with Debian 16.0-4-amd64 |
| POCASA | http://altair.sci.hokudai.ac.jp/g6/service/pocasa/ | 2010 | Available only as webserver |
| PocketAnalyzer | http://sourceforge.net/projects/papca/ or http://cpclab.uni-duesseldorf.de/downloads | 2011 | Not compatible with Debian 16.0-4-amd64 |
| QSiteFinder | http://www.bioinformatics.leeds.ac.uk/qsitefinder | 2005 | Not compatible with Debian 16.0-4-amd64 |
| Rate4Site | bioinfo.tau.ac.il/ConSurf | 2002 | Only available in Windows OS |
| SITEHOUND | http://scbx.mssm.edu/sitehound/sitehound-web/Input.html | 2009 | Not compatible with Debian 16.0-4-amd64 |
| SURFNET | http://www.ebi.ac.uk/thornton-srv/software/SURFNET/ | 1995 | Not compatible with Debian 16.0-4-amd64 |

 Table 5.2 Names, URLs, year of creation and respective limitations of the 18 pocket identification tools.

VMD, Visual Molecular Dynamics; PDB, protein data bank and OS, operating system.

| Group | Label | Description |
|---------------------------|-------------------------------|--|
| Ligand descriptors | lig_cov | Percentage of ligand covered by the predicted pocket |
| | poc_cov | Percentage of the pocket covered by the co-crystallized ligand |
| Size and shape | volume | Pocket volume in A^3 calculated via grid points |
| descriptors | surface | Pocket surface in A ² calculated via grid points |
| | lipo_surf | solvent accessible lipophilic surface; |
| | depth | Depth of the pocket in A |
| | ellips c/a or b/a | Ellipsoid main axes ratios, with $a > b > c$ |
| | enclosure | Rratio of number of surface to hull grid points |
| Functional group | H-don | Number of hydrogen bond donors |
| accomptore | H-acc | Number of hydrogen bond acceptors |
| | Met | Number of metals |
| | Hphob | Number of hydrophobic contacts |
| | siac ratio | Relative number of hydrophobic site interaction centers (SIACs, from flex) |
| Element descriptors | nof_dif_atms | Number of surface atoms lining the pocket |
| | elem_x | Number of elements of specific type in active site; types: C, N, O, S or other (X) |
| Amino acid composition | aa_apol, aa_pol and aa_neg | Relative number of amino acids apolar, polar, positive, and negative) |
| Amino acid | ALA, ARG, ASN, ASP, CYS, GLN, | Number of amino acids in pocket, 3-letter code of 20 amino acid |
| descriptors | GLU, GLY, HIS, ILE, LEU, LYS, | types |
| ' | MET, PHE, PRO, SER, THR, TRP, | |
| | TYR and VAL | |

 Table 5.3 List of the descriptors obtained from the DoGSiteScorer tool.

Table 5.4 List of the weight of each tool for each protein.

| | | GHECOM | Fnockot | DoGSiteScoror | lsoMif | Pro Δ C T 2 |
|--------------------|------------------|------------------|------------------|---------------|--------|-------------|
| | 0 /122 | | 0.2019 | 0.0677 | | 0.0362 |
| 1001 1001 | 0.4133 | 0.1901 0.2000 | 0.2010 | 0.0077 | 0.0071 | 0.0303 |
| | 0.4149 | 0.2000 | 0.1093 | 0.0711 | 0.00/0 | 0.0303 |
| 1 A D U 2 V V I | 0.3330 0.3330 | 0.1003 | 0.1420 0.1720 | 0.0323 | 0.1099 | 0.0303 |
| 211/0 | 0.3901 | 0.2005 | 0.1720 | 0.0007 | 0.1021 | 0.0327 |
| 2162 | 0.2049 | 0.0929 | 0.1101 | 0.1324 | 0.3243 | 0.0701 |
| 2002 | 0.2034 | 0.1420 | 0.2079 | 0.1029 | 0.1004 | 0.0507 |
| | 0.2709 | 0.1103 | 0.1900 | 0.1527 | 0.1975 | 0.0010 |
| | 0.2799 | 0.1049 | 0.1400 | 0.1010 | 0.2037 | 0.0640 |
| 1V3U | 0.20/3 | 0.1421 | 0.2110 | 0.1431 | 0.1710 | 0.0019 |
| 2000 | 0.3790 | 0.2107 | 0.1432 | 0.0767 | 0.1420 | 0.0453 |
| 30300 | 0.3021 | 0.1279 | 0.1229 | 0.1004 | 0.2303 | 0.0000 |
| 1003N | U.ZÕJ/ 0 2122 | U.1204 0.1222 | 0.2001 | 0.1017 | 0.1009 | 0.0032 |
| 200J 207T | 0.3433 | 0.133Z | 0.1343 | 0.1027 | 0.2234 | |
| 30Z1 | U.2024 | 0.1000 | U.ZZOJ 0 1022 | 0.1449 | 0.1003 | 0.0010 |
| 3663 | 0.3393 | 0.2004 0.1011 | U.1033 | | 0.1359 | 0.0331 |
| 2791 2055 | 0.3/94 | 0.1041 | 0.2701 | | 0.0743 | 0.0203 |
| | U.3443 | 0.2091 | 0.1097 | 0.0097 | 0.1499 | 0.03/3 |
| | U.2050 | 0.1090 | 0.1891 | 0.1047 | 0.10/5 | 0.0533 |
| | U.20/0 | | 0.1007 | 0.1399 | 0.2425 | 0.0763 |
| JE | 0.3170 | 0.0000 | 0.1060 | 0.1394 | 0.20/5 | 0.0010 |
| 313B | U.20UD | 0.1602 | 0.2000 | 0.1412 | 0.1333 | 0.0448 |
| 3004 | 0.3595 | 0.1491 | 0.1/15 | 0.0998 | 0.1794 | 0.0407 |
| 325X | 0.3/9/ | 0.1957 | 0.1857 | 0.0835 | 0.1116 | 0.0439 |
| | 0.3591 | 0.1501 | 0.1335 | 0.0890 | 0.2123 | 0.0561 |
| 3L4VV | 0.2539 | 0.1482 | 0.2218 | 0.1698 | 0.1647 | 0.0416 |
| | 0.2545 | 0.1071 | 0.1713 | 0.1459 | 0.2411 | 0.0801 |
| 32XF | 0.3847 | 0.1979 | 0.1989 | 0.0812 | 0.1083 | 0.0290 |
| | 0.4139 | 0.2012 | 0.1885 | 0.0714 | 0.0963 | 0.0287 |
| 30V1 | 0.2/98 | 0.1208 | 0.1/15 | 0.1646 | 0.21/6 | 0.0457 |
| 3HUU | 0.3/62 | 0.2115 | 0.1956 | 0.0845 | 0.0991 | 0.0331 |
| 20010 | 0.3581 | 0.1694 | 0.1975 | 0.0967 | 0.1398 | 0.0385 |
| 3F3C | 0.2674 | 0.1/53 | 0.2301 | 0.1415 | 0.13/2 | 0.0485 |
| 2ZXD | 0.3690 | 0.2211 | 0.1714 | 0.0840 | 0.1248 | 0.0297 |
| 2VVN | 0.2577 | 0.1213 | 0.2214 | 0.1546 | 0.1911 | 0.0539 |
| 3F3A | 0.3711 | 0.1865 | 0.1614 | 0.0977 | 0.1388 | 0.0445 |
| 2D3U | 0.3539 | 0.1333 | 0.1462 | 0.0964 | 0.2297 | 0.0405 |
| 1H23 | 0.3562 | 0.2298 | 0.1538 | 0.0793 | 0.1429 | 0.0381 |
| 2ZWZ | 0.3785 | 0.1522 | 0.1206 | 0.0819 | 0.2225 | 0.0442 |
| 3SU2 | 0.3553 | 0.1255 | 0.1467 | 0.1062 | 0.2203 | 0.0460 |
| 3BFU | 0.2563 | 0.1304 | 0.2397 | 0.1555 | 0.1566 | 0.0615 |
| 2P4Y | 0.4201 | 0.1800 | 0.2041 | 0.0750 | 0.0897 | 0.0311 |
| 1YC1 | 0.3806 | 0.1821 | 0.1683 | 0.0905 | 0.1289 | 0.0496 |
| Continued | | | | | | |

| 2J78 | 0.3884 | 0.2140 | 0.1579 | 0.0844 | 0.1135 | 0.0418 |
|------|--------|--------|--------|--------|--------|--------|
| 2X8Z | 0.3873 | 0.1739 | 0.2635 | 0.0687 | 0.0769 | 0.0297 |
| 3GBB | 0.4069 | 0.1970 | 0.1688 | 0.0785 | 0.1123 | 0.0365 |
| 2VW5 | 0.3799 | 0.1929 | 0.1595 | 0.0866 | 0.1300 | 0.0511 |
| 3GE7 | 0.3484 | 0.1782 | 0.2081 | 0.0867 | 0.1327 | 0.0458 |
| 1R5Y | 0.2391 | 0.1301 | 0.2643 | 0.1558 | 0.1491 | 0.0617 |
| 4DE2 | 0.2651 | 0.1292 | 0.2076 | 0.1613 | 0.1668 | 0.0701 |
| 3UDH | 0.4061 | 0.1951 | 0.2027 | 0.0788 | 0.0873 | 0.0299 |
| 3G2Z | 0.2581 | 0.1122 | 0.2129 | 0.1620 | 0.1860 | 0.0688 |
| 3EBP | 0.3483 | 0.1386 | 0.1527 | 0.1007 | 0.2143 | 0.0453 |
| 2GSS | 0.2590 | 0.1345 | 0.2499 | 0.1361 | 0.1651 | 0.0553 |
| 2XHM | 0.3878 | 0.1977 | 0.2380 | 0.0734 | 0.0809 | 0.0222 |
| 3MSS | 0.3987 | 0.2165 | 0.1761 | 0.0735 | 0.1015 | 0.0337 |
| 1U33 | 0.3343 | 0.1734 | 0.1813 | 0.1213 | 0.1471 | 0.0427 |
| 1SLN | 0.3597 | 0.1971 | 0.1919 | 0.0851 | 0.1304 | 0.0359 |
| 3F17 | 0.3821 | 0.1596 | 0.1650 | 0.1002 | 0.1493 | 0.0439 |
| 2X00 | 0.2487 | 0.1189 | 0.1581 | 0.1612 | 0.2579 | 0.0551 |
| 2BRB | 0.3511 | 0.2136 | 0.1689 | 0.0864 | 0.1398 | 0.0402 |
| 3AO4 | 0.2528 | 0.1397 | 0.3007 | 0.1416 | 0.1215 | 0.0437 |
| 4DE1 | 0.2556 | 0.1245 | 0.2134 | 0.1675 | 0.1739 | 0.0652 |
| 4GID | 0.4089 | 0.1918 | 0.2181 | 0.0707 | 0.0828 | 0.0277 |
| 3D4Z | 0.2929 | 0.0859 | 0.1403 | 0.1359 | 0.2788 | 0.0662 |
| 2ZCR | 0.4160 | 0.1872 | 0.1901 | 0.0751 | 0.0962 | 0.0354 |
| 3IMC | 0.3658 | 0.1946 | 0.1752 | 0.0920 | 0.1290 | 0.0434 |
| 10GS | 0.2425 | 0.1485 | 0.2463 | 0.1468 | 0.1689 | 0.0470 |
| 1GPK | 0.3622 | 0.2275 | 0.1535 | 0.0772 | 0.1415 | 0.0381 |
| 1PS3 | 0.3089 | 0.0839 | 0.1357 | 0.1309 | 0.2764 | 0.0643 |
| 3PE2 | 0.3659 | 0.1701 | 0.1457 | 0.0823 | 0.1911 | 0.0449 |
| 3ACW | 0.4033 | 0.1907 | 0.2083 | 0.0743 | 0.0889 | 0.0345 |
| 3FCQ | 0.2641 | 0.1451 | 0.2276 | 0.1387 | 0.1711 | 0.0535 |
| 3CYX | 0.3889 | 0.1417 | 0.1471 | 0.0891 | 0.1770 | 0.0563 |
| 2CET | 0.3820 | 0.2205 | 0.1569 | 0.0859 | 0.1138 | 0.0409 |
| 3K5V | 0.3689 | 0.2124 | 0.1660 | 0.0861 | 0.1275 | 0.0392 |
| 2V00 | 0.3695 | 0.1891 | 0.1833 | 0.0811 | 0.1287 | 0.0483 |
| 10S0 | 0.3317 | 0.1728 | 0.1963 | 0.0973 | 0.1623 | 0.0397 |
| 1Z95 | 0.4029 | 0.1514 | 0.1545 | 0.0835 | 0.1675 | 0.0403 |
| 3BKK | 0.3863 | 0.1892 | 0.2387 | 0.0697 | 0.0835 | 0.0327 |
| 3MFV | 0.3560 | 0.1737 | 0.1672 | 0.0955 | 0.1601 | 0.0474 |
| 2ZJW | 0.3655 | 0.2112 | 0.1799 | 0.0777 | 0.1320 | 0.0337 |
| 3DD0 | 0.2764 | 0.1277 | 0.1889 | 0.1473 | 0.2013 | 0.0585 |
| 3N86 | 0.2669 | 0.1487 | 0.2331 | 0.1452 | 0.1492 | 0.0569 |
| 3EHY | 0.3535 | 0.1721 | 0.1775 | 0.1085 | 0.1447 | 0.0436 |
| 3E93 | 0.3951 | 0.2141 | 0.1760 | 0.0691 | 0.1096 | 0.0361 |

| 1LOQ | 0.3578 | 0.2029 | 0.1504 | 0.1009 | 0.1505 | 0.0375 |
|------|--------|--------|--------|--------|--------|--------|
| 3F80 | 0.3698 | 0.1705 | 0.1628 | 0.0973 | 0.1491 | 0.0506 |
| 1U1B | 0.2616 | 0.1033 | 0.2079 | 0.1703 | 0.2027 | 0.0541 |
| 3NW9 | 0.2511 | 0.1711 | 0.2099 | 0.1463 | 0.1625 | 0.0591 |
| 3IVG | 0.4155 | 0.1983 | 0.1835 | 0.0734 | 0.0945 | 0.0347 |
| 3N7A | 0.2713 | 0.1502 | 0.2387 | 0.1393 | 0.1464 | 0.0541 |
| 1NVQ | 0.3624 | 0.1580 | 0.1569 | 0.0853 | 0.1903 | 0.0472 |
| 20BF | 0.3819 | 0.1906 | 0.1879 | 0.0873 | 0.1146 | 0.0377 |
| 3G0W | 0.3903 | 0.1891 | 0.1840 | 0.0842 | 0.1185 | 0.0339 |
| 1N2V | 0.3244 | 0.1948 | 0.1935 | 0.0965 | 0.1465 | 0.0443 |
| 3VH9 | 0.2863 | 0.1119 | 0.1313 | 0.1361 | 0.2704 | 0.0639 |
| 1LBK | 0.3877 | 0.2125 | 0.1763 | 0.0713 | 0.1189 | 0.0332 |
| 3AG9 | 0.3723 | 0.1294 | 0.1280 | 0.1019 | 0.2185 | 0.0499 |
| 4G8M | 0.3749 | 0.1991 | 0.1951 | 0.0865 | 0.1081 | 0.0361 |
| 2G70 | 0.3613 | 0.1873 | 0.1955 | 0.0948 | 0.1221 | 0.0390 |
| 3SU5 | 0.3655 | 0.1380 | 0.1343 | 0.1049 | 0.2136 | 0.0437 |
| 2VO5 | 0.2533 | 0.1543 | 0.2135 | 0.1504 | 0.1642 | 0.0643 |
| 3S8O | 0.2515 | 0.1237 | 0.1910 | 0.1645 | 0.2228 | 0.0465 |
| 3B68 | 0.3807 | 0.1694 | 0.1623 | 0.0819 | 0.1648 | 0.0410 |
| 3PWW | 0.3497 | 0.1793 | 0.1529 | 0.1021 | 0.1746 | 0.0414 |
| 2XY9 | 0.3882 | 0.1905 | 0.2272 | 0.0776 | 0.0839 | 0.0325 |
| 3ZSO | 0.3119 | 0.1780 | 0.2135 | 0.1119 | 0.1237 | 0.0609 |
| 2FVD | 0.3752 | 0.2055 | 0.2073 | 0.0765 | 0.1058 | 0.0297 |
| 1KEL | 0.3009 | 0.1216 | 0.2132 | 0.1412 | 0.1685 | 0.0545 |
| 3UEX | 0.3623 | 0.1668 | 0.1639 | 0.0943 | 0.1785 | 0.0342 |
| 2ZCQ | 0.4172 | 0.1855 | 0.1895 | 0.0747 | 0.1019 | 0.0311 |
| 2YGE | 0.3661 | 0.1964 | 0.1641 | 0.0809 | 0.1414 | 0.0511 |
| 3SU3 | 0.3632 | 0.1373 | 0.1341 | 0.1067 | 0.2134 | 0.0453 |
| 1SQA | 0.2332 | 0.1385 | 0.2138 | 0.1915 | 0.1757 | 0.0472 |
| 2D10 | 0.3658 | 0.1863 | 0.1786 | 0.0900 | 0.1384 | 0.0409 |
| 2VOT | 0.4032 | 0.1359 | 0.1175 | 0.0817 | 0.2115 | 0.0503 |
| 2WEG | 0.4151 | 0.1397 | 0.1293 | 0.0887 | 0.1809 | 0.0463 |
| 2XNB | 0.3553 | 0.1969 | 0.1881 | 0.0851 | 0.1378 | 0.0368 |
| 3B3S | 0.2495 | 0.0982 | 0.1517 | 0.1665 | 0.2637 | 0.0704 |
| 3CJ2 | 0.2575 | 0.0997 | 0.1736 | 0.1493 | 0.2597 | 0.0602 |
| 3FK1 | 0.2670 | 0.1479 | 0.2725 | 0.1439 | 0.1109 | 0.0578 |
| 4TMN | 0.2402 | 0.1521 | 0.2409 | 0.1543 | 0.1514 | 0.0610 |
| 4DJR | 0.3558 | 0.1428 | 0.1231 | 0.1099 | 0.2175 | 0.0509 |
| 1P1Q | 0.4020 | 0.1787 | 0.2029 | 0.0812 | 0.0987 | 0.0365 |
| 2XDL | 0.3686 | 0.1537 | 0.1479 | 0.0972 | 0.1790 | 0.0537 |
| 4GQQ | 0.2517 | 0.0979 | 0.1211 | 0.1647 | 0.2953 | 0.0695 |
| 3GNW | 0.3769 | 0.2130 | 0.1803 | 0.0793 | 0.1181 | 0.0324 |
| 3DXG | 0.2593 | 0.1136 | 0.2285 | 0.1737 | 0.1713 | 0.0535 |

| 2X0Y | 0.2764 | 0.1535 | 0.2337 | 0.1458 | 0.1389 | 0.0516 |
|------|--------|--------|--------|--------|--------|--------|
| 3MUZ | 0.2579 | 0.1597 | 0.2509 | 0.1475 | 0.1367 | 0.0473 |
| 1HNN | 0.3899 | 0.1849 | 0.1960 | 0.0812 | 0.1089 | 0.0391 |
| 3EJR | 0.2715 | 0.0863 | 0.1446 | 0.1441 | 0.2895 | 0.0640 |
| 2QFT | 0.4253 | 0.1908 | 0.1888 | 0.0744 | 0.0817 | 0.0389 |
| 2YMD | 0.2671 | 0.1069 | 0.1763 | 0.1738 | 0.2201 | 0.0557 |
| 2QMJ | 0.3813 | 0.1194 | 0.1027 | 0.0948 | 0.2598 | 0.0420 |
| 2XYS | 0.3611 | 0.1461 | 0.1326 | 0.0991 | 0.2165 | 0.0446 |
| 3VD4 | 0.2621 | 0.1286 | 0.1520 | 0.1450 | 0.2527 | 0.0596 |
| 2HB1 | 0.2612 | 0.1244 | 0.2735 | 0.1507 | 0.1431 | 0.0471 |
| 1HFS | 0.2732 | 0.1491 | 0.2359 | 0.1328 | 0.1661 | 0.0430 |
| 3L3N | 0.3594 | 0.2314 | 0.1736 | 0.0763 | 0.1234 | 0.0359 |
| 30WJ | 0.3834 | 0.1963 | 0.1951 | 0.0727 | 0.1128 | 0.0397 |
| 3L7B | 0.3913 | 0.1925 | 0.2401 | 0.0686 | 0.0771 | 0.0304 |
| 1W3L | 0.2803 | 0.1223 | 0.1995 | 0.1619 | 0.1733 | 0.0627 |
| 3COY | 0.4210 | 0.1946 | 0.1849 | 0.0743 | 0.0918 | 0.0333 |
| 3NQ3 | 0.3653 | 0.1797 | 0.1717 | 0.0968 | 0.1520 | 0.0345 |
| 3NOX | 0.2513 | 0.0961 | 0.1297 | 0.1398 | 0.3114 | 0.0717 |
| 2ZX6 | 0.3739 | 0.1797 | 0.1630 | 0.0853 | 0.1582 | 0.0399 |
| 3LKA | 0.3635 | 0.1766 | 0.1665 | 0.1055 | 0.1439 | 0.0441 |
| 3FV1 | 0.3933 | 0.2019 | 0.1709 | 0.0787 | 0.1163 | 0.0388 |
| 3MYG | 0.3527 | 0.2009 | 0.1989 | 0.0861 | 0.1263 | 0.0351 |
| 2PQ9 | 0.4107 | 0.1839 | 0.2146 | 0.0721 | 0.0837 | 0.0351 |
| 3KV2 | 0.3541 | 0.1728 | 0.1762 | 0.0965 | 0.1539 | 0.0464 |
| 20LE | 0.2768 | 0.1041 | 0.1787 | 0.1509 | 0.2360 | 0.0535 |
| 1LOL | 0.3349 | 0.1709 | 0.1585 | 0.1129 | 0.1865 | 0.0363 |
| 1JYQ | 0.2774 | 0.0991 | 0.1475 | 0.1526 | 0.2646 | 0.0589 |
| 2IWX | 0.4143 | 0.2165 | 0.1744 | 0.0714 | 0.0933 | 0.0301 |
| 3L4U | 0.2695 | 0.0859 | 0.1381 | 0.1432 | 0.3048 | 0.0585 |
| 2VL4 | 0.2695 | 0.1279 | 0.1900 | 0.1509 | 0.1963 | 0.0653 |
| 3KWA | 0.2726 | 0.1229 | 0.1911 | 0.1413 | 0.2076 | 0.0645 |
| 4DJV | 0.2601 | 0.1085 | 0.1834 | 0.1655 | 0.2207 | 0.0617 |
| 3U9Q | 0.4046 | 0.2014 | 0.2074 | 0.0717 | 0.0879 | 0.0270 |
| 3UEU | 0.3611 | 0.1628 | 0.1580 | 0.0961 | 0.1789 | 0.0430 |
| 2WCA | 0.3901 | 0.1985 | 0.1615 | 0.0863 | 0.1305 | 0.0331 |
| 1E66 | 0.3604 | 0.2369 | 0.1500 | 0.0743 | 0.1390 | 0.0394 |
| 2XB8 | 0.2739 | 0.1603 | 0.2293 | 0.1365 | 0.1443 | 0.0557 |
| 3G2N | 0.3911 | 0.1987 | 0.2305 | 0.0737 | 0.0744 | 0.0317 |
| 1QI0 | 0.2703 | 0.1301 | 0.2058 | 0.1737 | 0.1650 | 0.0551 |