

# Protein structure modeling for variant pathogenicity prediction

Author: Sylt Schuurmans

Studentnumber: 333332

Study: Bio-informatica

Hanze hogeschool: Institute for Life Science and Technology

Supervisor: Martijn Herber

Hospital: University Medical Center Groningen: Department of Genetics

Supervisor : Joeri van der Velde

June 12, 2019

## **Abstract**

Around 1 in 17 people is affected by one of 7,000 known rare diseases. Most of these patients do not receive a diagnosis, which means they remain in uncertainty without a prognosis, are unable to join specific patient support groups, and do not receive the most appropriate treatment. Next-generation sequencing (NGS) of DNA promises to establish a molecular diagnosis and help these patients but many challenges still stand in the way of maximum success. Recent years have seen great advances in computational tools that quickly reduce the amount of DNA variants to be interpreted by a human expert for potentially pathogenic effects [1]. Although algorithms can now safely remove around 95% of the harmless variants, this still leaves hundreds of variants to be investigated for a whole-exome sequenced patient, which is far too many for a quick and clear diagnosis. Current tools to predict variant pathogenicity rely on features such as evolutionary conservation, annotation of regulatory genomics elements or structural DNA features. These tools have already been optimized over many years and further significant improvements are not expected. Therefore there is still a great need for even more powerful variant prioritization tools. A refreshing alternative was presented by VIPUR [2] which shows the potential of structural modeling of proteins to predict the actual effect of a specific variant on the function of that protein. This presents an exciting new opportunity to improve genome diagnostic variant prioritization. However, this predictor was (i) not integrated with the latest and greatest variant pathogenicity prediction approaches, (ii) was trained on relatively small number of variants, and (iii) did not result in high quality software that was ready to be taken into routine diagnostic practice. To test this approach we will explore the potential pitfalls of protein modeling by evaluating the VIPUR pipeline and by examining a single protein with its variants.

## **Acknowledgment**

This report has been written for the genomic coordination center (GCC) to gain insight into structural data to improve the GAVIN variant predictor. I want to thank Joeri van der Velde for supervising me during this project and helping me with writing this report. I also would like to thank Benjamin Kant and Marielle van Gijn for helping me decide to focus on assessing tumor necrosis factor associated receptor-associated periodic syndrome (TRAPS). And I especially would like to thank Tsjerk Wassenaar for informing us about the function and structure of TNFRSF1A, giving us the appropriate protein structures to work with and steering this project into a meaningful direction.

## Abbreviations

3D Three Dimensional  
ACCP Solvent Accessible Surface Area  
API Application Programming Interface  
Bash Bourne Again Shell  
CPU Central Processing Unit  
CSV Comma Separated Values  
DNA Deoxyribonucleic Acid  
FADD Fas Associated Death Domain protein  
FEM Fixed End Move  
FHF Familial Hibernian Fever  
GAVIN Gene-Aware Variant INterpretation  
GCC Genomic Coordination Center  
GRCh/hg Genome Reference Consortium Human Human Genome  
HOPE Have yOur Protein Explained  
LOMETS Local Meta-threading Server  
MD Molecular Dynamics  
MPI Message Parsing Interface  
NCBI National Center for Biotechnology Information  
NF- $\kappa$ B Nuclear Factor kappa-light-chain-enhancer of activated B cells  
OpenGL Open Graphics Library  
OS Operating System  
OSF Open Science Framework  
PDB Protein Data Bank  
PM Pivot Movement  
PSI-BLAST Position Specific Iterative BLAST  
PSSM Position Specific Scoring Matrix  
RCSB Research Collaboratory for Structural Bioinformatics  
REU Rosetta Energy Unit  
RNA Ribonucleic acid  
RMSD Root Mean Square Deviation  
SASA Solvent Accessible Surface Area  
SCOP The Structural Classification of Proteins  
SLURM Simple Linux Utility for Resource Management  
SODD Silencer of Death Domain  
SPVAA Simple Protein Variant Analysis Approach  
TNF Tumor Necrosis Factor  
TNF $\alpha$  Tumor Necrosis Factor Alpha  
TNF $\alpha$  Tumor Necrosis Factor Alpha  
TNF $\beta$  Tumor Necrosis Factor Beta  
TNFB Tumor Necrosis Factor Beta  
TNFR1 Tumor Necrosis Factor Receptor Superfamily Member 1A TNFRSF1A Tumor Necrosis Factor Receptor Superfamily Member 1A  
TRADD Tumor Necrosis Factor Receptor type 1-Associated DEATH Domain protein  
TRAPS Tumor necrosis factor associated Receptor-Associated Periodic Syndrome  
VIPUR Variant Interpretation Using Rosetta  
VTS VIPUR Training Set

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Mutations and its effects in the central dogma of molecular biology . . . . .	1
1.2	A general concept of structural levels within proteins and the effect of mutations . . . . .	1
1.3	Addition of structural data to diagnosis and treatment in healthcare . . . . .	1
<b>2</b>	<b>Theory</b>	<b>2</b>
2.1	Protein modeling techniques . . . . .	2
2.2	A theoretical large scale implementation of structural protein variant assessment . . . . .	2
2.3	Monte Carlo method . . . . .	3
2.4	The use of the Monte Carlo method and its pitfalls . . . . .	3
2.5	Tumor Necrosis Factor Receptor Associated Syndrome . . . . .	4
2.6	Tumor Necrosis Factor Receptor Super Family Member 1A . . . . .	4
2.7	Tumor Necrosis Factor Alpha and Beta . . . . .	4
<b>3</b>	<b>Materials and methods</b>	<b>5</b>
3.1	VIPUR approach . . . . .	5
3.2	Rosetta . . . . .	5
3.2.1	Relax . . . . .	5
3.2.2	DDG Monomer . . . . .	6
3.2.3	Rescore . . . . .	6
3.2.4	Backrub . . . . .	6
3.3	PyRosetta . . . . .	6
3.4	PSI-BLAST . . . . .	6
3.5	Probe . . . . .	6
3.6	Robetta prediction server . . . . .	7
3.7	I-TASSER prediction server . . . . .	7
3.8	Modeller . . . . .	7
3.9	GAVIN Machine Learning Data Table . . . . .	7
3.10	GenomAD . . . . .	7
3.11	Infevers . . . . .	7
3.12	Research Collaboratory for Structural Bioinformatics . . . . .	7
3.13	Uniprot . . . . .	8
3.14	PyMOL . . . . .	8
3.15	HOPE . . . . .	8
3.16	Bash . . . . .	8
3.17	Python . . . . .	8
3.18	R scripting language . . . . .	9
3.19	SLURM . . . . .	9
3.20	MPI . . . . .	9
<b>4</b>	<b>Results</b>	<b>10</b>
4.1	Reviving the VIPUR approach to expand rare disease diagnostics . . . . .	10
4.1.1	Preparatory steps for using the VIPUR approach . . . . .	10
4.1.2	Resolving VIPUR system incompatibilities . . . . .	10
4.1.3	Expanding the VIPUR training set with data from TNFRSF1A by homology modeling and protein threading . . . . .	11
4.1.4	Practical VIPUR usage . . . . .	11
4.2	Analyses of proteins variants TNFRSF1A . . . . .	12
4.2.1	Requirements for determining structural and binding effects of protein variants . .	12
4.2.2	Introduction of the simple protein variant analysis approach . . . . .	12
4.2.3	Carrying out SPVAA on TNFRSF1A . . . . .	13

4.3	Finding mutation information with HOPE . . . . .	19
<b>5</b>	<b>Discussion</b>	<b>20</b>
<b>6</b>	<b>Conclusion</b>	<b>22</b>
<b>7</b>	<b>Future work</b>	<b>22</b>

## List of Figures

1	Flowcharts VIPUR pipeline and altered VIPUR pipeline . . . . .	10
2	I-TASSER and Robetta models with and without templates . . . . .	11
3	Flowchart SPVAA pipeline . . . . .	13
4	TNFRSF1A homotrimer with TNF $\alpha$ homotrimer relax density plots . . . . .	15
5	TNFRSF1A homotrimer with TNF $\beta$ homotrimer relax density plots . . . . .	16
6	TNFRSF1A homotrimer with TNF $\alpha$ homo trimers wild type and mutated relaxed models	17
7	TNFRSF1A homotrimer with TNF $\beta$ homo trimers wild type and mutated relaxed models	18

## **List of Tables**

1	Sample from the combined observed TNFRSF1A mutations table . . . . .	13
2	Sample of the TNFRSF1A PDB residue mutation table . . . . .	14

# **1 Introduction: Variant prediction in genome diagnostics and the addition of protein modeling**

## **1.1 Mutations and its effects in the central dogma of molecular biology**

Within the human genome mutations occur continuously by internal and external factors that: insert, remove, substitute or alter the reading frame in a nucleotide sequence. Mutations are not without consequences and can be protective [3], benign or harmful by altering the deoxyribonucleic acid (DNA) order. From a sequence of DNA genes are transcribed into ribonucleic acid (RNA) which can work as machinery or translates into an amino acid sequence to form a protein. Mutations outside a gene could lead to lowered or heightened transcription of a protein, when a mutation resides inside a gene it could lead to proteins that are unstable during or after formation, perform less optimal or are not functional [4–6].

## **1.2 A general concept of structural levels within proteins and the effect of mutations**

The formation of protein structures is classified in different levels, distinctions are made based on bindings and structures that arise with the interaction of bonds. The order in which amino acids appear in a sequence is called the primary structure, in this level amino acids are only bound to each other by peptide bonds. Within a primary structure amino acids can form new peptide bonds between the N-terminus and C-terminus of an amino acid, with these bonds 3D structures are made called  $\alpha$ -helices and  $\beta$ -sheets that together make up the secondary structure. The tertiary structure gives further rise to the 3D shape of a polypeptide by making disulfide bridges, ion and hydrogen -bonds, hydrophobic and hydrophilic -interactions between amino acids By combining multiple tertiary structures the quaternary structure of a protein can be formed out of the mentioned bonds, bridges and interactions [7, 8].

Mutations within proteins can have different effects to protein structures, often single missense mutations often have minimal effect on the backbone of a protein [9, 10] but can result in destabilization of the structure when assembled or can disrupt the active site. Frameshift mutations on the other hand can cause large differences in the primary structure and have therefore a higher chance of an altered sequence that leads to deformation or stop codon introduction [11].

## **1.3 Addition of structural data to diagnosis and treatment in healthcare**

Acquiring information about DNA sequences depends on sequencing, which became cheaper over the years [12], and found its use in diagnosing patients within the healthcare sector [1]. From the collected data by genome sequencing experiments most of the analysis is handled in-silico due to the quantities of data that are produced [1]. Proteins often find their use in diagnosing diseases experimentally [8, 13], however in-silico it is often limited to information about conservation in the amino acid sequence [14]. Yet, the 3D shape of proteins defines their function [15] and by assessing structures it can become possible to determine changes in function that are caused by mutations that might not be discoverable through conservation and are therefore unclassifiable. Another advantage of structural information is that it becomes possible to develop treatment with diagnostic information for diseases that are caused by mutations [16]. With experimental methods such as X-ray crystallography and nuclear magnetic resonance (NMR) more than 158000 structures [17] have been completely revealed, however it is only a tiny fraction of the potential possible proteins [18] (especially without the inclusion of all folds). Making 3D structures is currently not common for diagnosis because it is relative expensive and is difficult to perform, some structures contain flexible regions which makes it hard to determine the exact position of some atoms and can cause information loss about the structure [19, 20].

## 2 Theory

### 2.1 Protein modeling techniques

An alternative approach for determining structures, compared to experimentally is

An alternative approach to determine structures is based on modeling protein structures computationally from the amino acid sequence . A downside from computer generated models is that they do not follow the laws of physics and therefore not automatically fold into the correct confirmation.

With the method homology modeling sequences of the requested protein are aligned to sequences of known experimental determined structures, based on these alignments a template is formed whereon structural fragments are built, it is not recommended to use this strategy if the sequence identity is less than 20% since there might not be any structural relation at that point [10]. Another approach is protein threading which relies on the observation of folds in previous determined experimental structures. Based on the occurrence of specific folds a probability is predicted that a certain residue in a protein might fold in that manner.

Strategies are continuously being improved and developed for proteins to determine the unknown structures, but all have the similar guidelines in avoiding steric hindrance [21] and finding the lowest energies based on different scoring systems [22]. From the computer generated models many are less accurate than the experimental determined methods and are often compared to them for reference. However the computational models do not have follow the experimental laws of physics which bottleneck the current experimental methods in throughput, but also for example in producing structures that represent membrane proteins [23].

### 2.2 A theoretical large scale implementation of structural protein variant assessment

With the wide spectrum of potential different proteins it can be difficult and maybe momentarily impossible to produce any form of universal protein assessment standardization that is able to determine if a mutation is harmful or not based on structural information. However a first step to solve such a complex problem would be by determining the correct approach, in this case it is assumed that a machine learning approach would be the best method for detecting patterns in structures and classifying the effect of structural changes. Because it has the ability to learn from structural mutations currently available, assuming that the current knowledge about structures and mutations is correct, and is able to develop new insights in how structural changes could affect proteins.

Since the problem is so complex it should be divided into smaller more feasible problems, beginning by separating the different protein classes, which for example can be done according to The Structural Classification of Proteins database (SCOP) [24]. A first discrimination between the proteins could be made based on protein type/fold class (membrane, globular, fibrous and disordered -proteins) because these differences already predetermine some functions and locations for certain proteins in a cell [25–28]. After formation of these classes each should have its own machine learning method applied so their features can be analyzed within context of where and how they function. The next set of discriminators is highly dependent on the variations in classes, but all have features in the end describing bonds, interactions and movement of complexes in protein structures. When for each of the main classes a method has been developed a meta classifier will determine based on certain aspects which method should be applied to determine the effect of mutation in a protein.

### **2.3 Monte Carlo method**

There are complex problems in a variety of research fields which could take up years or even centuries to compute with simple deterministic methods. For some problems there is an algorithm which makes it possible to cut down computation time significantly, but when no deterministic algorithm is available to speed up the process an empirical probabilistic method might be able to approximate the desired result. With the Monte Carlo method random samples are taken from the parameter space ,that describe a data set, and fed into a model which produces a potential outcome. By repeating the process more results are generated until at some point the data can display a pattern that describes the outcome. The result is a quantified probability which describes the chance that something might occur based on the quantity of occurrence generated by the model [29–31].

The Monte Carlo methods can differ depending on the algorithm and application in which it is used, but in summary most implementations will follow a general pattern [30]:

0. Construct a model which is able to describe an outcome of the problem.
1. Define the space of which inputs can be used by the model to get an outcome (creating a parameter space).
2. Use the model to generate results based on random sampled input from the parameter space.
3. Order and determine which results are part of a certain outcome and draw conclusions on the generated statistical evidence.

### **2.4 The use of the Monte Carlo method and its pitfalls**

The Monte Carlo method is widely used within various applications in different fields of science but it is limited in the type of problems it can solve and is suitable for; problems of which all the inputs are known but it is too inefficient to compute deterministically; situations that require uncertainty to be incorporated into the analysis and exploring parameters for a model that give a better impact than the current parameters. The mentioned type of problems it can solve all tend to rely on significant quantities of data which makes it a relative time consuming process for generating results. Meaning of the generated result is highly depended on the model and random sampling techniques which both contribute to an errors in the result [29, 30, 32].

## **2.5 Tumor Necrosis Factor Receptor Associated Syndrome**

Tumor necrosis factor receptor-associated periodic syndrome (TRAPS) is classified as a rare disease (1 : 1,000,000) and was formerly known as Familial Hibernian fever (FHF) [33], is a hereditary autosomal dominant disease which can cause recurring fevers with a duration from days up to several months. Symptoms during these fevers are: skin rash, swelling, inflammatory reactions across the whole body and pain in the abdomen, muscles and/or joints, a long term and lasting effect is the accumulation of amyloid within the kidneys and may result in other diseases [6]. TRAPS is known to be caused by mutations within the gene tumor necrosis factor receptor 1 (TNFRSF1A/TNRF1) (Section 2.6), the mutated proteins tend to get trapped in the cell and will be unable to reach the cell surface and therefore start activating a inflammatory response [6, 34]. So far 158 mutations have been associated with the disease [35], but more mutations have been identified in TNFRSF1A wherein some might be pathogenic (Sections 3.9, 3.10).

## **2.6 Tumor Necrosis Factor Receptor Super Family Member 1A**

Tumor Necrosis Factor Receptor Super Family Member 1A (TNFRSF1A, TNFR1) is a gene located on chromosome 12 region 1 band 3 and sub-band 31. The gene produces a trans-membrane receptor consisting of 445 residues divided into 221 residue cytoplasmic section and a 171 extracellular part that consists of 4 conserved cysteine rich domains [36–38]. The receptor is ubiquitous across most cell surfaces ,but not on erythrocytes [39], and can form two different types of unbound hexagonal clusters depending on the dimer formation [40]. When the structures are dimers the binding sites are exposed and make it possible for tumor necrosis factor (TNF)  $\alpha$  and  $\beta$  (Section 2.7) to bind in trimeric form, with binding of TNF the dimers disconnect and three TNFR1s interact with the TNF trimer [40]. With the interaction of the TNF trimers with TNFR1 it can activate several pathways such as; the nuclear factor kappa-light-chain-enhancer of activated B cells (NF- $\kappa$ B), which enhances the transcription of various genes during inflammation, infection or other forms of external stress; also it is able to activate the extrinsic pathway of apoptosis after binding of TNF to TNFR1, by releasing the silencer of death domain (SODD) proteins release on the cytoplasmic site. Tumor Necrosis Factor Receptor type 1-Associated DEATH Domain protein (TRADD) [41] will start to bind together with proteins that will form a complex which will attract Fas associated death domain (FADD) and after two hours[42] if not inhibited. On binding of FADD initiator caspase 8 starts a cascade wherein caspase 3 is activated an will cleave aspartate out of proteins and thereby disrupting the metabolism[43–45].

## **2.7 Tumor Necrosis Factor Alpha and Beta**

The proteins TNF  $\alpha$  and  $\beta$  are both pro-inflammatory cytokines that are produced as response to an infection or when a cell is damaged. Both are transcribed from their genes that reside in chromosome 6 in the p-arm at region 2 band 1 and sub-band 3. TNF  $\alpha$  and  $\beta$  are 35% identical and 50% homologous to each other consisting out of 233 and 205 amino acid residues. Both are able to form a homotrimeric structures that can bind to the dimeric form TNFR1 (Section 2.6) to activate the extrinsic pathway[37, 46–48].

## 3 Materials and methods

### 3.1 VIPUR approach

Variant interpretation using Rosetta (VIPUR) is a machine learning approach for predicting deleteriousness of proteins (loss of function) and uses sequential and structural information. To train VIPUR a training set was made that contains sequence and structure features. Structures within the VIPUR training set (VTS) were collected from Modbase [49] and SWISS-MODEL [50–54]. Proteins that did not have a structure within the VTS were modeled with Modeller (Section 3.8 based on protein fragments that had the highest amino acid sequence identity to the protein. Some structures from the databases had: duplicate chains, ligands, metals and non-standard amino acids which were removed to avoid inconsistencies that could alter the features generated by the tools and hinder feature collection within Rosetta tools (Section 3.2), therefore they were removed to make the data homogeneous. For all proteins a mutation file was made that described where and which residues had to be mutated, with this file DDG monomer (Section 3.2.2) can predict changes in the protein structure of which all features are used for prediction. Structural mutations of proteins that are in the VTS were introduced by a script using PyMOL (Section 3.14) by default or PyRosetta (Section 3.3) if PyMOL was not available. After mutation the structure is optimized by the relax application (Section 3.2.1) and produces 50 relaxed structures of a single variant where the properties of each protein are written to a score file of which the quartiles are used as a learning feature. Probe (Section 3.5) calculated the solvent accessible surface area (SASA/ACCP) of a protein in square Ångstrom ( $\text{\AA}^2$ ) which is a structural machine learning feature. The sequence features of VIPUR are produced by PSI-blast (Section 3.4) on non mutated sequences and blasted against the NCBI protein database (nr) which results in a position specific scoring matrix (PSSM). From the PSSM scores of the non-mutated, mutated, the difference in scores between non-mutated and mutated, information content and the difference between groups [2, 55] are all sequential features for VIPUR. With 106 features generated by the mentioned tools deleteriousness of a protein variant is determined with sparse logistic regression. The term sparse implies that a limited set of features was used because the weights "shrink" to 0 with regularization [56].

### 3.2 Rosetta

Rosetta is a software suite that has a variety of tools that are developed to aid in macro molecular and antibody analysis, design and prediction [57]. However no tools in the suite have been encountered that could introduce missense mutations in the proteins and has been dealt with by other software (Sections 3.14, 3.3, 3.8). With the introduced mutations water had to be removed because some tools cannot predict structures well with: water, metals and amino acids that are no part of the standard (20) amino acids [58].

Within the tools from Rosetta various scores are assigned to different properties related to bonds, interactions, energies and geometries within structures and are written to a score file. From all different scoring metrics the Rosetta score or Rosetta energy unit (REU) , which is the total\_score in the score files, can be used to compare models made from the same protein per tool. Not only is the score based on energy but also it has statistical terms which influence the score based on known favorable folds from existing structures that reside in the curated Rosetta database [23]. In summary a lower Rosetta score would make a more natural model.

*Rosetta software suite Version 3.10*

#### 3.2.1 Relax

The Relax application was used by VIPUR and by SPVAA to relax the side chains to minimize energy levels within the local conformational search space [59] of the structure, it determines the most likely energy levels with the Monte Carlo method (Section 2.4) and after a certain set of moves it produces a structure and starts anew [60, 61]. Scores from each produced by Relax were written into a single score file.

### 3.2.2 DDG Monomer

DDG monomer is meant to predict energetic stability of a point mutation in monomeric protein. The application was used by VIPUR to collect features related to energies and hydrogen, disulfide, bonds and constraints differences between the wild type and a mutated protein. To execute the tool a script had to be ran that renames the wild type pdb file and it requires a "mutation file" that describes the change of a residue based on name and position changes to a different residue [62].

### 3.2.3 Rescore

With this tool Rosetta scores can be calculated based on silent or PDB files proteins structures [63] , the output is identical to that is written within the score files produced by Relax (Section 3.2.1).

### 3.2.4 Backrub

The backrub application is based on the Monte Carlo method (Section 2.4), and alters a protein by moving its backbone residues with a strategy called fix end move (FEM). With this strategy, groups of residues are selected at random from the structure, it can contain up to: four dihedral, two bond angles and two end points. Both ends of a group are fixated at their position in which a new angle  $\alpha$  arises, within this angle residues are pivoted in their natural occurring maximum range of  $\pm 10^\circ$  [64, 65]. With this method the backbones of newly introduced mutations were altered, for each attempt a new file was generated and a score was written to a score file, from which the lowest Rosetta scoring was selected to be further relaxed (Section 3.2.1). It was used on the mutated protein to relax the modified backbone structure.

## 3.3 PyRosetta

Is an application programming (API) which has Python bindings (Section 3.17) for the Rosetta software suite (Section 3.2), it founds its use in VIPUR when no PyMOL (Section 3.14) was available to mutate residues within a structure [66].

*Version 4*

## 3.4 PSI-BLAST

Position specific iterative basic local alignment search tool (PSI-BLAST) focuses on distant relatives of proteins by making a profile of the sequence and querying it at a protein sequence database. With the generated results a new profile is constructed and queried again, these steps are repeated several times to determine which residues are found in relatives of the protein. The result is a position specific scoring matrix (PSSM) which describes the frequency of which residues are substituted by a specific other residue, positive is more, negative is less common [67–69]. From the PSSMs sequences features were acquired for the VIPUR machine learning method.

*Position-Specific Iterated BLAST 2.7.1+*

## 3.5 Probe

Probe is able to evaluate atom packing for a single protein or interacting proteins by creating a probe, which is described as a sphere like object, that marks an area with dots when at least two non-covalent atoms are in contact with the probe at the same position [70, 71]. VIPUR used this tool to calculate solvent accessible surface area (SASA or ACCP).

*version 2.16.130520*

### **3.6 Robetta prediction server**

The web tool Robetta integrates several tools to form protein structures based on sequence alignments of previously discovered structures also known as homology modeling (Section 2.1). It requires an amino acid sequence, optionally constrains and fragments can be added to disallow movement of certain structures or add known fragments to avoid calculating pieces that are already known. With this information Robetta search with the help of sequence aligners for known fragments and tries to incorporate them into a single protein structure [72–76].

### **3.7 I-TASSER prediction server**

The I-TASSER web server is a tool that is able to predict protein structures with a FASTA sequence. The first step it takes is finding structural templates which resemble the sequence by local meta-threading server (LOMETS). LOMETS starts with multiple sequence alignment of which several sequences will undergo protein threading by different programs to form structural templates. The templates are assessed based on the highest alignment Z-score, the program specific confidence score and sequence identity [77, 78]. The known fragments of TNFRSF1A (Section 2.1) were given as a template to I-TASSER and modeled into a whole protein to make it possible to introduce mutations and predict pathogenicity of a variants.

*Server version*

### **3.8 Modeller**

Modeller is software that is developed for homology modeling but it was used for its utilities which allowed to; complete protein data bank (PDB) structures with missing atoms; predict disulfide bonds that were missing and mutate protein residues [79, 80].

*Version 9.21*

### **3.9 GAVIN Machine Learning Data Table**

Is a collection of nucleotide mutations from rare diseases used by the GAVIN [1] machine learning approach. From this set the genes of TNFRSF1A (Section 2.7) with a missense mutation were filtered (Section 3.18) and written into a format which the variant effect predictor could (VEP) [81] could read and translate from nucleotide to protein mutations. The classification of these variants was according to Clinvar significance values [82].

### **3.10 GenomAD**

The GenomAD database consists of unified data from large scale genome sequencing data projects and is based on genome reference consortium human genome build 37 human genome 19 (GRCh37/hg19). From this database missense mutations were collected for TNFRSF1A (Section 2.6), no classification was known for these mutations [83].

### **3.11 Infevers**

Is a website about hereditary auto immune diseases with for each disease a downloadable table about the known mutations and their classification, when classified with pathogenic or benign its function has been observed. The table for TRAPS disease (Section 2.6) was used to collect missense mutations of TNFRSF1A gene [84].

### **3.12 Research Collaboratory for Structural Bioinformatics**

Research Collaboratory for Structural Bioinformatics (RCSB) is a database where whole or fragmented experimentally determined proteins structures that are published can be found and downloaded. The

Fragments for modeling (Sections 3.6, 3.7) whole TNFRSF1A (Section 2.6) (1EXT [85]) and determining the differences in energy levels (Section 3.2.1) with TNF  $\beta$  (1TNR [38]) with the interaction site were acquired from this database [86].

### 3.13 Uniprot

Knowledge from various omic domains about proteins has been linked together into single database called Uniprot which makes all information accessible at once, for TNFRSF1A (Section 2.7) the FASTA sequences were collected from Uniprot and for structures it redirected to (Section 3.12) [87].

### 3.14 PyMOL

Visualization of 3D structures; making images of proteins; putting the known orientations of monomeres in position; replacing TNF  $\beta$  structure with a TNF  $\alpha$  in the bound structure and aligning the structures to measure the distance between X-ray crystal structures and the produced models were done with PyMOL [88]. PyMOL had a different use in VIPUR where it was used in combination with Python (Section 3.17) to perform mutagenesis on the protein structures to introduce a missense mutations.

*Version 2.2.3*

### 3.15 HOPE

Have yOur Protein Explained (HOPE) is a web service that collects information of about a user specified missense mutation in a protein and comes from various sources. Uniprot (Section 3.13) is queried with BLAST to find homologous sequences and structures, other features that are found on Uniprot are active sites, domains and various other sequence features that help to identify the function of a region. From the BLAST results homology models are made with Yasara that are sent of to WHAT IF web services that calculate structural information about the protein. Before the formation of a report all information is put into a decision tree to asses mutational effects in contexts of: contacts, structural locations, non-structural features, previous variant information and amino acid properties to form an automated report [89–92]. With this method it is not possible to asses ligands and complexes at once but only a single missense mutation within a monomer. *Version 1.1.1*

### 3.16 Bash

Unix like operating systems (OS) have a shell which allows users to interact with programs on a computer or with the computer itself based on commands submitted. The default shell for MacOS and also for several Linux distributions is the Bourne again shell (Bash) which was used to launch Python scripts (Section 3.17) and submit jobs to the SLURM workload manager (Section 3.19).

*Laptop Version GNU bash, version 3.2.57(1)-release (x86\_64-apple-darwin18)*

*Server Version GNU bash, version 4.1.2(2)-release (x86\_64-redhat-linux-gnu)*

### 3.17 Python

Both VIPUR and the pipeline that minimizes backbone (Section 3.2.4) and side chain energies (Section 3.2.1) were written in Python due to its capabilities, ease of use and because modeller (Section 3.8) for Mac OS relies on the system version of Python and does currently not support newer versions besides the one found within the OS of Mac. The mutations that were put together from the different tables (Sections 3.10, 3.11, 3.9) with R Section 3.18) were filtered by a Python script. To apply each mutation correctly on the proteins in the detailed method a script was written in which files were generated that described in a compact format on which chains and position a mutation resided.

*Laptop version 2.7.15*

*Server version 2.7.11*

### **3.18 R scripting language**

With R the tables from GenomeAD, GAVIN and Infevers (Sections 3.10 3.9 3.11) of TNFRSF1A missense mutations (Section 2.6) were merged together in a new comma seperated values file with their known classifications. Ordering and filtering the double mutations and removing double classifications where done with Python (Section 3.17). It has also been used in combination ggplot2 [93] and data.table [94] to make density plots of all scores acquired from Rosetta Backrub and Relax.

*R scripting front-end version 3.5.2 (2018-12-20)*

### **3.19 SLURM**

For computational jobs where a laptop or desktop does not suffice because due to the lack computational resources a computer cluster could come to aid. These clusters consist out of several computers that execute resource intensive tasks, to manage these systems for many clients and to use these clusters optimal a workload manager mlike simple Linux utility resource management (SLURM), is installed. Jobs are submitted that request resources for execution and are scheduled on the systems queue which is ordered based on priorities, resource requirement and time.

### **3.20 MPI**

Some tools from the Rosetta software suite (Sections 3.2) have the ability to use multiple central processing unit (CPU) cores from a single computer or from multiple computers. With a message parsing interface (MPI) it is possible for software to communicate between CPU cores on the same and on different computers to exchange information about processes and therefor solving solutions faster.

*OpenMPI/1.8.8-GNU-4.9.3-2.25*

## 4 Results

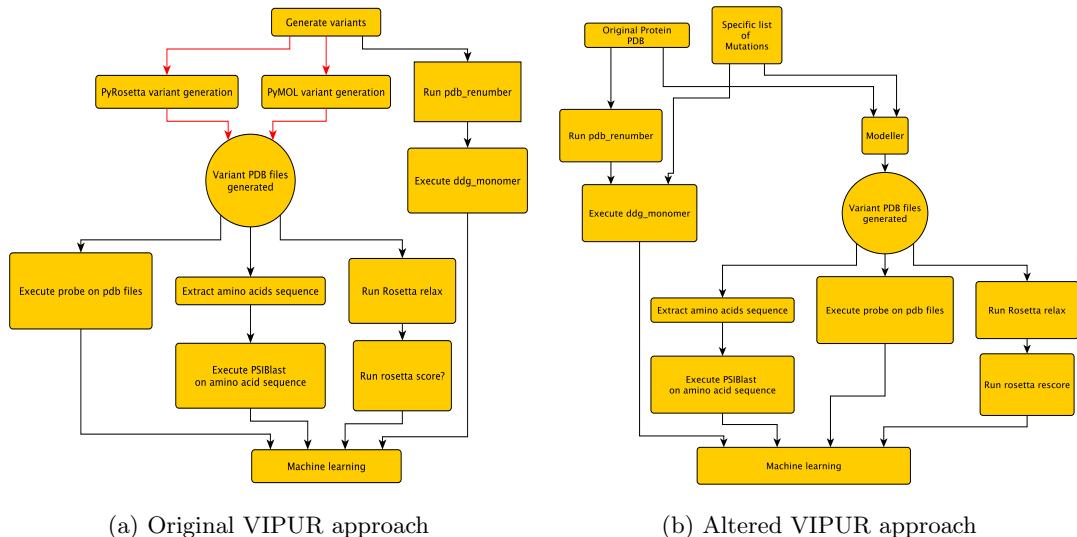
### 4.1 Reviving the VIPUR approach to expand rare disease diagnostics

#### 4.1.1 Preparatory steps for using the VIPUR approach

After the publication of VIPUR the tools, data and applications became available at the open science framework (OSF) [95] which were downloaded and reviewed. All applications from the Rosetta software suite (Section 3.2) were pre-compiled without support for MPI (Section 3.20) and with that not the ability to benefit from multiple CPUs. The Rosetta software suite was rebuilt with MPI support in a slurm job where the compilation could benefit from multiple CPU cores.

#### 4.1.2 Resolving VIPUR system incompatibilities

Within the VIPUR pipeline residues were mutated to determine the effects of a structural mutation, by default missense mutations were inserted with PyMOL (Section 3.14), an alternative method integrated within the pipeline for situations wherein PyMOL was not accessible Pyrosetta (Section 3.3) could be used. Neither of these programs could be built or compiled because the lack of Open graphics library (OpenGL) for PyMOL and having the incorrect C++ and C libraries for PyRosetta. To bypass both programs and still be able to introduce mutations into PDB files Modeller (Secton 3.8) was introduced and built.



(a) Original VIPUR approach

(b) Altered VIPUR approach

Figure 1: Both flowcharts illustrate the VIPUR pipeline wherein each block is a procedure the central circle is the purpose of the mutated applications and each arrow represents the path to it. Figure 1a has red arrows that indicate that both methods were incapable to produce the mutated PDB files. Within figure 1b the alternative method is proposed wherein PyMOL and PyRosetta (Sections 3.14, 3.3) is substituted by Modeller (Section 3.8) to acquire the mutated protein structures.(To zoom in on the details within the figure it is recommend to look at the PDF version.)

#### 4.1.3 Expanding the VIPUR training set with data from TNFRSF1A by homology modeling and protein threading

Since the VTS did not have any features of TNFRSF1A (Section 2.6) the amino acid sequence was collected from Uniprot (Section 3.13) and the protein from RCSB (Section 3.12). The structures available of TNFRSF1A were incomplete, fragments for the TNF  $\alpha$  and  $\beta$  binding site [38] were available and its death domain that interacts with TRADD [96] which plays a role in apoptosis (Section 2.6). To acquire a monomeric structure of TNFRSF1A two ab initio modeling web services I-TASSER and Robetta (Sections 3.7, 3.6) had been employed. Both were given the task to model the whole protein with and without a template to determine how well they could model a known structure and what it would form. Determination of the best model was based on the smallest root mean square deviation distance (RMSD) in Å, between a produced model compared to the X-ray crystallographic model of the TNFRSF1A binding site.

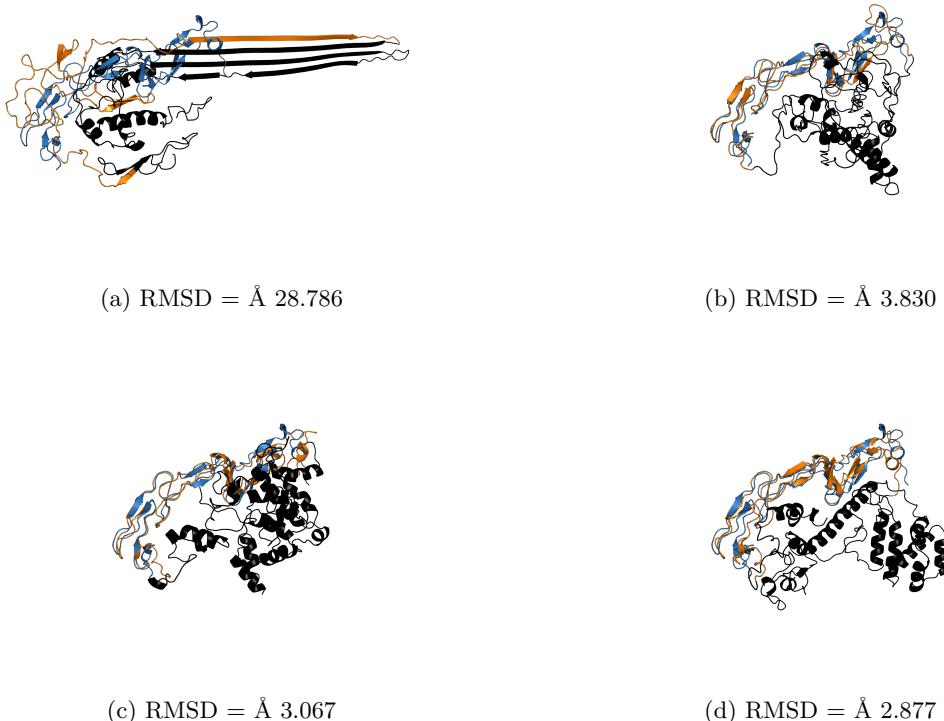


Figure 2: 3D structures of TNFRSF1A ( 2a, 2b: I-TASSER, 2c, 2d: Robetta) without (left: 2a, 2c) and with templates (right: 2b, 2d). The sky blue colored structure in each figure is an X-ray crystallographic model (1EXT) of the binding site of TNFRSF1A and the orange structure is the representation of that identical fragment in the model made by the web services.(To zoom in on the details within the figure it is recommend to look at the PDF version.)

#### 4.1.4 Practical VIPUR usage

With the hindrance of software incompatibility on the cluster, difference in produced models between the web services, discovery of consequences by removing elements from structures, and the time it would take to reverse engineer VIPUR a new decision was formed. VIPUR would be set aside for now and if time was left there would be further looked into to make it applicable for protein structure evaluation.

## 4.2 Analyses of proteins variants TNFRSF1A

### 4.2.1 Requirements for determining structural and binding effects of protein variants

Protein variants can be assessed from multiple perspectives and together they can form a holistic view on how a protein works and how mutations affect its workings. However adding perspectives to the protein assessment makes it complex and requires expertise to determine its validity and contribution, therefore the analysis has been limited to basic structural information and also make the assessment inline with the VIPUR methods.

Various proteins consist of multiple chains that can be identical or different depending on their function [97] and should be taken into account when assessing protein variants since one residue might alter the binding between chains and might alter the proteins formation. Different molecules and atoms that do not make up a protein but play a role in a pathway and function (ligands and co-receptors) are able to affect a proteins shape and can behave differently when a residue is mutated.

A different aspect that can change with mutations is the alteration in motions between structures which can allow or disallow certain movements to occur and inhibit processes.

### 4.2.2 Introduction of the simple protein variant analysis approach

A different method for determining function loss in a protein variant is through assessment of difference in energy levels between a wild type and a variant in the complex where it resides. Analyzing mutations from this perspective gives the ability to view a protein in whole and determine how residues cause perturbations in a protein. To make a variant of the wild type, a structure was required wherein a missense mutation could be introduced with Modeller (Section 3.8). The backbone structure of the variant was modified with the backrub application (Section 3.2.4) to make it better interact with other amino acid backbones in the structure resulting in 1000 models. The lowest Rosetta scoring (Section 3.2) one would be selected to further improve the side chains with the Relax application (Section 3.2.1) which makes the side chains of the protein move into lower energy formations. From these 64 models were made where its energies levels were compared to the native structure to determine the effects a mutation would have on the protein, the lowest scoring ones were used as figures. This method shows similarities to that of VIPUR and was only tested on TNFRSF1A (Section 2.6) and its ligands TNF  $\alpha$  and  $\beta$ . This method keeps: duplicate chains and protein ligands within the structure, water is excluded since it can cause issues with Rosetta tools (Section 3.2). This procedure with these steps is defined as the Simple protein variant analysis approach (SPVAA).

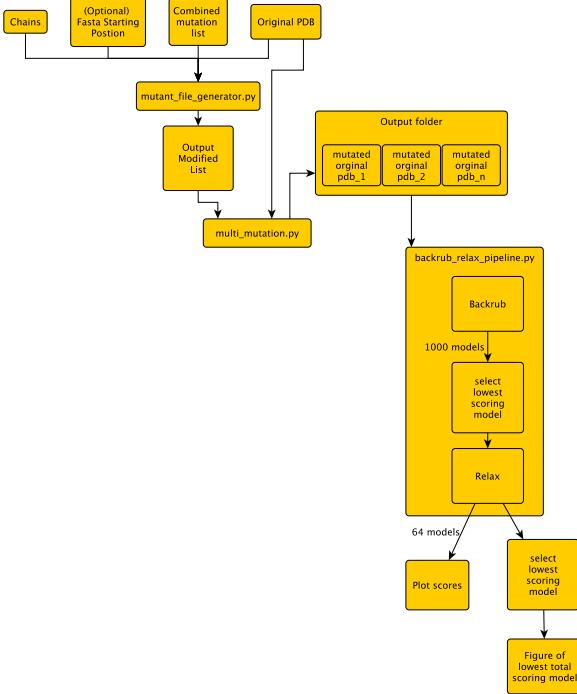


Figure 3: Flowchart of SPVAA wherein a list of known mutations generate the appropriate information for modeller to mutate residues in the original PDB of the protein and feed it into the backrub relax pipeline where the model is altered to go into a lower energy state.(To zoom in on the details within the figure it is recommend to look at the PDF version.)

#### 4.2.3 Carrying out SPVAA on TNFRSF1A

Before introducing mutations into a protein structure it is helpful to know if a mutation has been observed to avoid allocating resources to something that does not occur. Therefore three tables with observed TNFRSF1A mutations (Sections 3.9, 3.10, 3.11) have been combined with an R script (Section 3.18) into a single table consisting of two columns. The first column (split into three columns 1) contains strings that describes the: original residue, position and where it mutates to, the second column describes whether a formed mutation is harmful, with most mutations the effects have not been identified yet.

Original residue	Position in the protein sequence	New residue	Classification
Cys	44	Tyr	PATHOGENIC
Thr	44	Pro	PATHOGENIC
Thr	44	Ser	PATHOGENIC

Table 1: The format wherein mutations were filtered from the GAVIN, GenomAD and Infevers tables (Sections 3.9, 3.10, 3.11), describe whether a structural mutation is harmful or not. For many mutations it is unknown and other classifications are available, to view the whole table visit the supplementary.

For assessing variants in TNFRSF1A a structural fragment was used that contained TNF  $\beta$  (1TNR) [38] and was made homotrimeric with PyMOL (Section 3.14) which results in six chains that emulate a bound TNFRSF1A with TNF  $\beta$ . The first column of the mutation table did not contain sufficient information to apply mutations correctly and within the PDB different numbering is used than in the amino acid sequence. To bundle the information and make it usable for introducing mutations a Python script (Section 3.17) has been written that combines the mutation table, PDB chains and the correct

position within the sequence into a type of table which has sufficient information to mutate structures.

Iteration number	Filename	Chain	Residue index in chain	New residue
34	1tnr3_TNFA	R	0	TYR
34	1tnr3_TNFA	T	0	TYR
34	1tnr3_TNFA	S	0	TYR
35	1tnr3_TNFA	R	0	PRO
35	1tnr3_TNFA	T	0	PRO
35	1tnr3_TNFA	S	0	PRO
36	1tnr3_TNFA	R	0	SER
36	1tnr3_TNFA	T	0	SER
36	1tnr3_TNFA	S	0	SER

Table 2: The format that describes the mutations that should be made by Modeller (Section 3.8), with specifications of the: model, file, chain, residue index and the new residue. The whole table for TNFA and TNFB are visible within the supplementary.

To introduce mutations within PDB structures a Python script (Section 3.17) was written which used the generated mutation table (Table: 2) and a matching PDB structure, from the table; it acquires an iteration number which specifies if a mutation has to be stored in a single file or across multiple files; the filename serves as key which determine the PDB that should be used; letters specify chains, numbers are indices within the chains and the last column states the three letter residue where it should mutate to. When a structure is read in through the Python bindings of Modeller (Section 3.8) all non standard atoms and molecules are removed because Rosetta (Section 3.2) is not able to deal with those atoms. Just before mutagenesis takes place missing atoms are added to the structure that were difficult to determine with experimental methods(Section 1.3). After the insertion of all mutations a last attempt was made by modeller to add disulfide bridges, based on the distance between cysteine residues, into the structure.

With the many protein variants generated and limited resources to use SPVAA mutations had to be chosen. All mutations were picked from the Infevers table because that were variants of a single isoform and had a protein structure (1TNR [38]) available that interacted with TNF  $\beta$ . Mutations cysteine 62 to glycine and phenylalanine acid 141 to isoleucine were validated as pathogenic mutations within Infevers table and were used to determine the effectiveness of SPVAA. Within the infevers table no benign validated missense mutations were available, but to still have the opportunity to asses a likely benign mutation, the mutation glutamic acid 138 to alanine has been added to SPVAA.

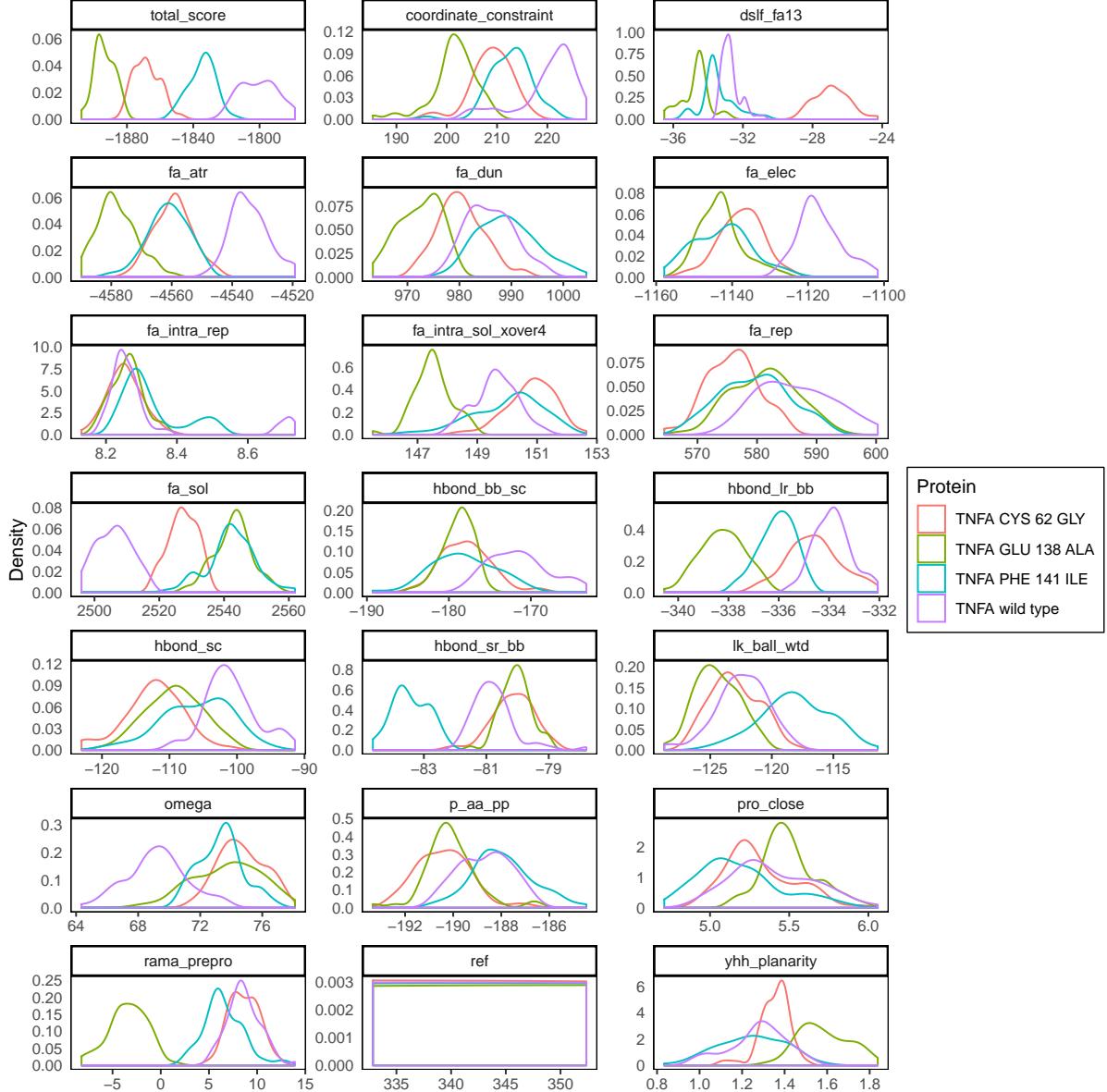


Figure 4: Density of different scoring metrics from the models produced with relax of the wild type and all mutants that interacted with TNF $\alpha$ . Most of the total scores, fa\_atr values and fa\_elec values of the mutated model are lower than the wild type that is bound to TNF $\alpha$  and the fa\_sol values are higher of the wild type than the mutants. TNFA CYS 62 GLY has more higher values at the dslf\_fa13 (disulfide geometry potential) than all other mutations. (Plots of backrub TNF $\alpha$  scores are in the supplementary.) (To zoom in on the details within the figure it is recommended to look at the PDF version.)

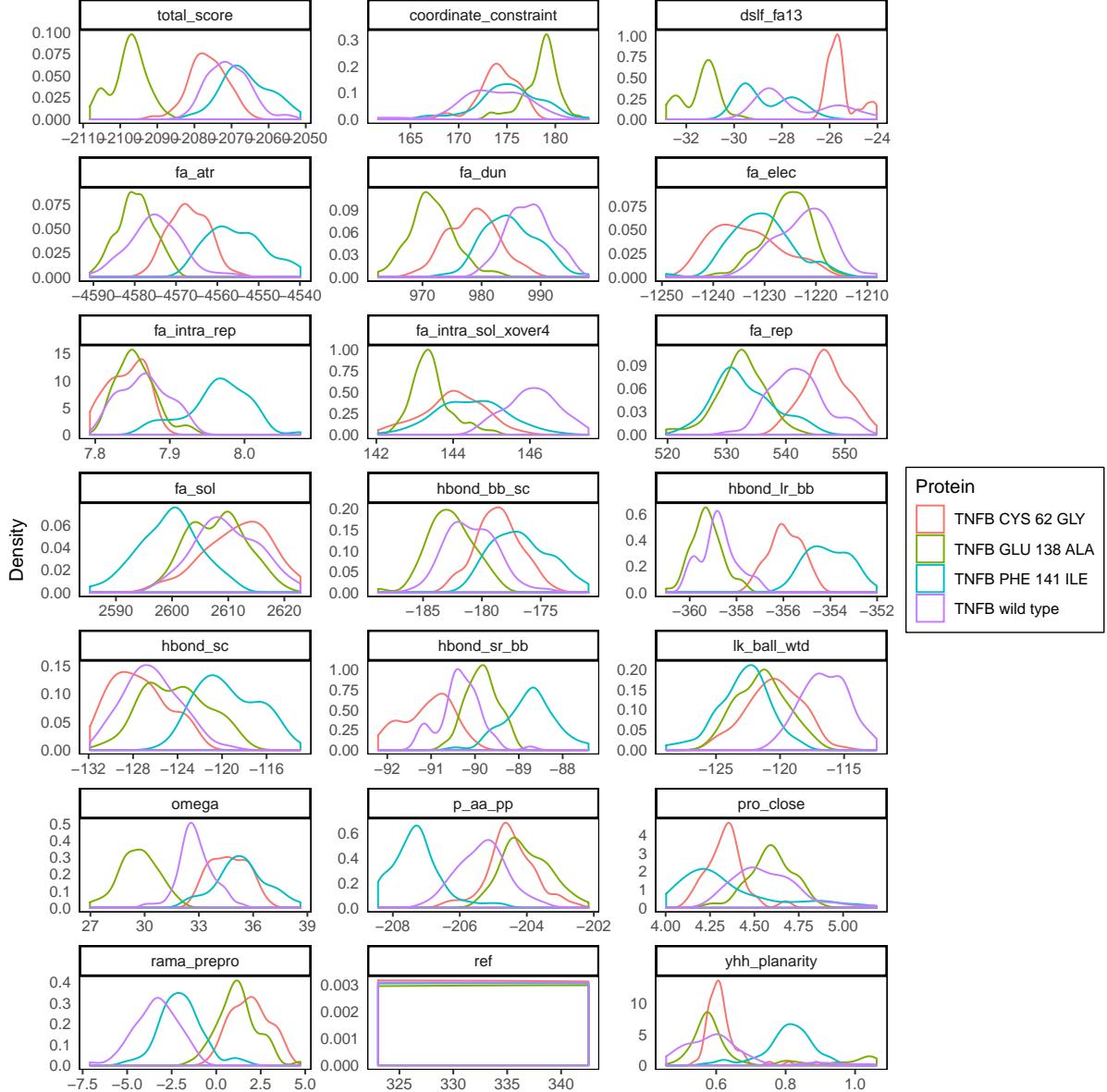


Figure 5: Density distributions of the scores generated by relax for the wild type and mutations of TNFRSF1A that interact with TNF $\beta$ . The total score values of TNFB GLU 138 ALA are lower than all other models and dslf\_fa13 (disulfide geometry potential) values are higher at TNFB CYS 62 GLY than the other models. fa\_intra\_rep (Lennard-Jones repulsive between atoms in the same residue) is higher within the models of PHE 141 ILE and hbond\_lr\_bb (Backbone-backbone hbonds distant in primary sequence) is higher at CYS 62 GLY and PHE 141 ILE. (Plots of backrub TNF $\beta$  scores are in the supplementary.) (To zoom in on the details within the figure it is recommended to look at the PDF version.)

In the attempt to make mutated structures behave more natural two tools from the Rosetta software suite (Section 3.2) had been used to minimize energies within protein structures. With the Backrub application (Section 3.2.4) 1000 altered backbone models have been produced each with 10000 Monte Carlo moves (Sections 2.4). For each model that Backrub generated a set of scores were assigned to the properties, which together formed a collective score that described energy and bond occurrence in nature (Section 3.2). Models of the wildtypes and mutants with the lowest collective score ,the best score, were chosen to undergo further side chain optimization within the Relax (Section 3.2.1). 64 different relaxed models were produced and with various scores related to the properties of which the one with lowest total score would be chosen to visualize.

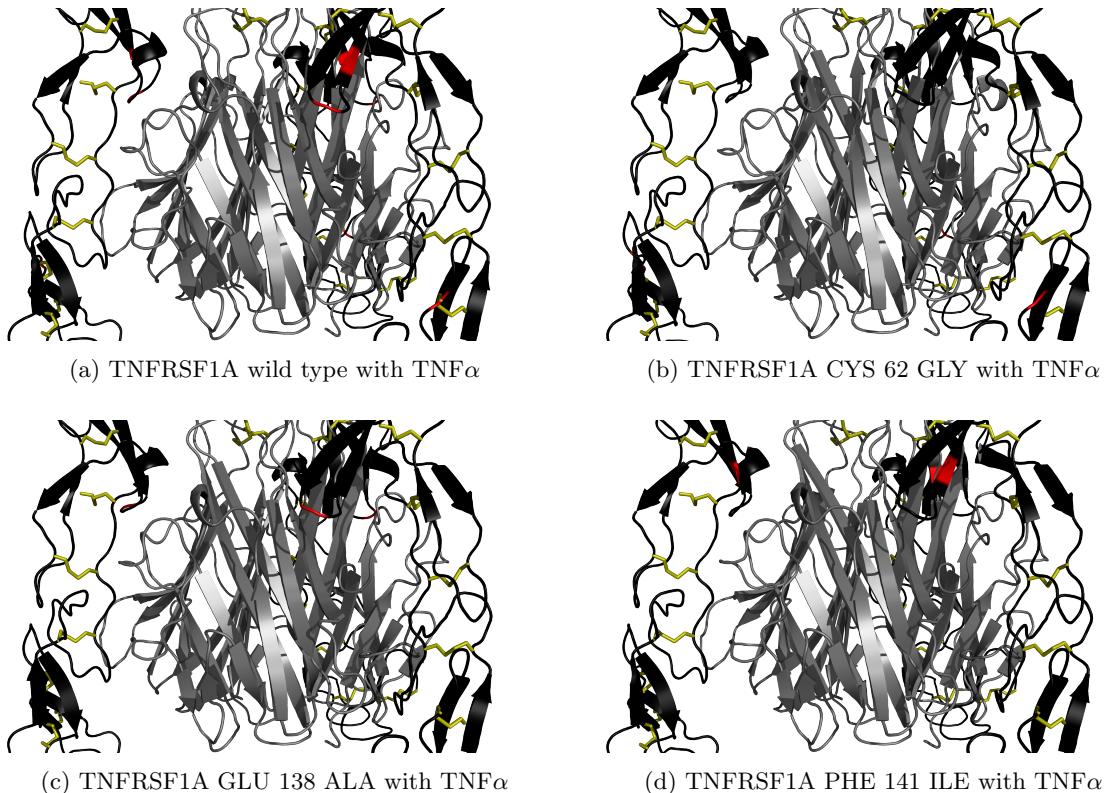


Figure 6: 3D structures of a homotrimer TNFRSF1As (black) with a homotrimer TNF $\alpha$ s (gray) and disulfide bridges (dark yellow). The wild type (6a) has three red colored areas which are the original residues of the protein before any form of mutation. Within CYS 62 GLY (6b) it is visible that at the position where a mutation is introduced (red) a disulfide bridge is missing. The mutations of GLU 138 ALA (6c) and PHE 141 ILE (6d) show no large differences at the mutated spots (red).(To zoom in on the details within the figure it is recommend to look at the PDF version.)

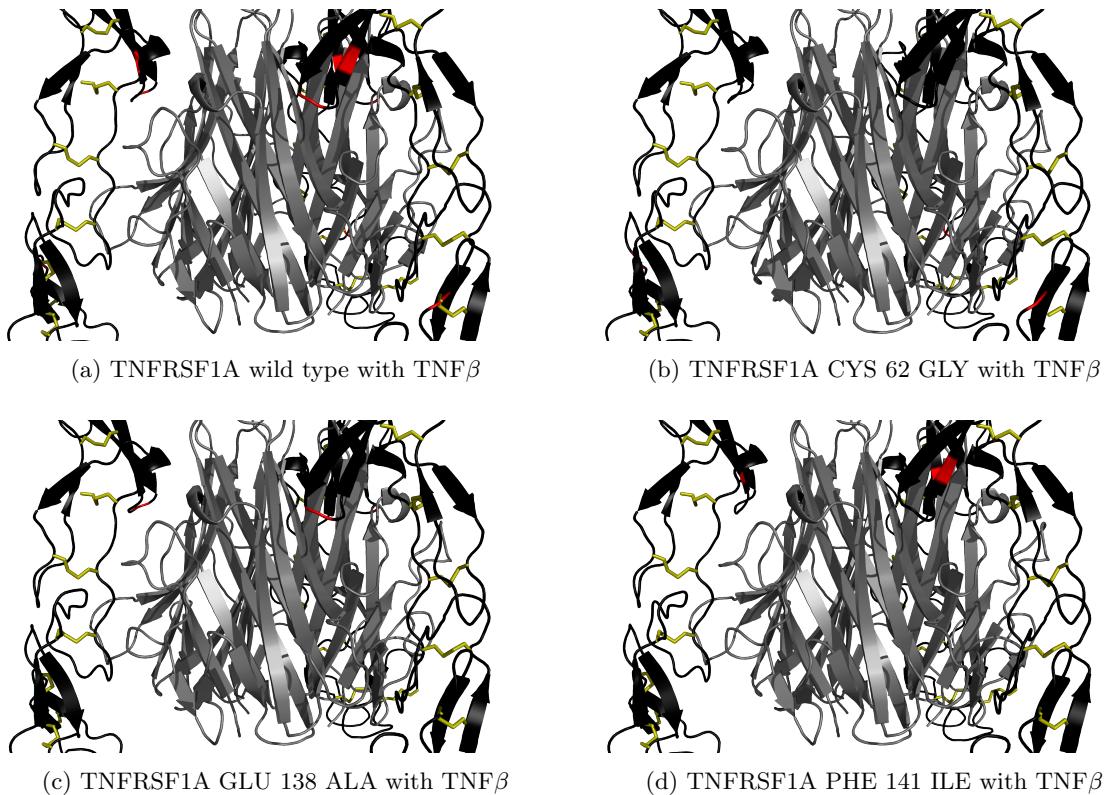


Figure 7: 3D structures of a homotrimer TNFRSF1As (black) with a homotrimer TNF $\beta$ s (gray) and disulfide bridges (dark yellow). The wild type (7a) has three red colored areas which are the original residues of the protein before any form of mutation. Within CYS 62 GLY (7b) it is visible that at the position where a mutation is introduced (red) a disulfide bridge is missing. The mutations of GLU 138 ALA (7c) and PHE 141 ILE (7d) show no large differences at the mutated spots (red). (To zoom in on the details within the figure it is recommend to look at the PDF version.)

### 4.3 Finding mutation information with HOPE

With uncertainty in the numbers provided by SPVAA a more textual informative approach was used with the web service HOPE (Section 3.15). The mutations that were known from Infevers (Section 3.11) and were also used with SPVAA were tested by HOPE (CYS62GLY, GLU138ALA and PHE141IIE), all reports are visible within the supplementary.

HOPEs first test was the mutation of Cysteine 62 to glycine, which is known within the Infevers table as pathogenic and was validated. It discovered that the residue was involved in a disulfide bridge and was 100% conserved in related protein sequences, based on the observation that cysteine formed a disulfide bridge it expected that with the replacement of it glycine would make the whole structures less rigid. HOPE predicts that mutation is pathogenic because of the high conservation of the residue, which is further confirmed by its search results in which it found the original publication where the observation has been described and associated to TRAPS [35].

According to Infevers is the mutation of glutamic acid at position 138 mutated to alanine classified as likely benign and was not validated yet. HOPE discovered with a BLAST query that glutamic acid occurs often at position but other residues such as alanine have been observed at the position. Structurally glutamic acid forms salt bridges with proline 368 and leucine 390 and is found in a sequence of amino acids that is repeated throughout TNFRSF1A. The amino acid lies within a domain where it interacts with other domains and is important for the protein's activity, with this mutation it might already perturb the binding capabilities according to HOPE.

The last mutation that was tested with HOPE was phenylalanine 141 to Isoleucine and was according to Infevers pathogenic and has been validated. Phenylalanine has been conserved at this position and few other residues have been seen at the position, it is a member of the identical domain as glutamic acid 138 and HOPE predicts that it would not damage the protein based on this information. However within the structure it could inhibit interaction with other domains and protein activity.

## 5 Discussion

People with rare diseases are currently hard to diagnose and are often not put in the desired treatment groups, with machine learning methods such as GAVIN 95% of the benign variants can be harmfully removed, however these methods rely on conservation and have been heavily optimized over the years [1]. A new refreshing approach called VIPUR uses sequential and structural data to predict deleteriousness of a protein variants.

Within the attempt to make VIPUR usable for the diagnosis of rare diseases it was discovered that some questionable steps were taken to make it especially applicable for diagnosis but also to determine deleteriousness of proteins itself: (i) "All protein models were standardized to remove unwanted components (duplicate chains, ligands, metals and non-standard amino acids)" [2]. Standardizing data can be beneficial to avoid learning features from proteins that are available to some models but should not be the determining factor for classification. However any form of context to the protein is removed and might therefore make incorrect assumptions about how: a monomer interacts with other monomers, ligands, metals, non-standard amino acids and water which can all have an effect on how proteins shape and interact [98]. (ii) With the utilization of Rosetta's Relax application different models are formed based on the Monte Carlo method (Section 2.4), VIPUR produces 50 structures with Relax per protein which is a tiny amount of the potential search space of possible folds that could have made changes in a mutated protein, which also can be seen in the scores of the model made from of TNFRSF1A with TNF  $\alpha$  &  $\beta$  (Figures 4, 5), Rosetta itself suggests to make sufficient models, starting with 5000[99]. (iii) The features acquired with probe in combination with the models that were produced, within the publication of Probe is mentioned: "It requires both highly accurate structures and also the explicit inclusion of all hydrogen atoms and their van der Waals interactions." [70]. It is currently not possible to determine if the structures were accurate, however is it likely that no loose hydrogen atoms were included within the structure based on the knowledge that all structures were standardized and likely some of the structures did not have any loose hydrogens within them. To make the outcome of Probe useful to VIPUR the program Reduce should have ran first, which add hydrogen atoms to the structure, which is recommend on the site were Probe can be downloaded from [71].

More questions arise when further investigating the publication; within the figures (4, 5) [2] and supplementary figures (10, 11) [56] are heatmaps of PSSMs added that display the values of the natural and mutated residues. In combination with the methods used on standardizing structures and collecting features a suspicion arises that there is little contribution from the structural features and that prediction depends on PSI-BLAST results. Figure 4 within the publication shows a protein wherein residue serine 204 which part of an  $\alpha$  helix is mutated to proline and is predicted as deleterioueness by VIPUR, which is logical even without predictions because prolines are known to be  $\alpha$  helix stoppers [100] and therefor affect the  $\alpha$  helix it form.

VIPUR has not been used for various reasons; The models that were predicted by the web services (Figures 2) had a decent accuracy for the binding site, but the rest of structures differ in various sections and make it hard to determine if these were accurate representations of TNFRSF1A; VIPUR was built on and for a single system and required reverse engineering to test the basic demos and PyMOL or PyRosetta (Sections 3.14, 3.3) still had to be replaced with modeller (Section 3.8) which would have required more time to repair. The differences in between the models are probably due to the fact that TNFRSF1A is transmembrane protein which is hard to acquire structures from with experimental methods [21], even though  $\sim 57\%$  of TNFRSF1As structure was known (1EXT[85] , 1ICH[96]), if more models would have been produced a potential accurate structure could have been formed .

Although it is not part of developing a technique that can help diagnosing rare disease variants but the publication contains a claim ("VIPUR can be applied to mutations in any organism's proteome...." [2]) which contradicts with its methods: "remove unwanted components (duplicate chains, ligands, metals and non-standard amino acids)" [2]. Currently there are more than 140 amino acids found in natural proteins of which 22 are part of the amino acid alphabet and 20 of those are classified as standard [101]. By removing the non-standard amino acids from proteins it becomes impossible to analyze mutations in every organisms proteome.

SPVAA did not assess whole complex and neither did became a machine learning tool ready to use in diagnosis to make automated predictions. The protein information have similar weaknesses as VIPUR wherein: water, metals and other molecules are removed from its structure. However the proteins can keep their extra monomers and protein ligands, even when the structure requires identical or different ones they can be added manually into the structure and analyzed.

For three variants it has been attempted to acquire structural information to determine their differences with the wild type structure, in two of the three structures hardly any changes were visible. From the models that were assessed only homotrimeric TNF $\alpha$  - $\beta$  bound mutations were processed, not the unboud dimeric structures. From the proteins that were modeled too few were produced of each structure, backrub used 10000 Monte Carlo moves, which is very little compared to the amount of residues it had and should have required more. The relax application made 64 models per mutation but Rosetta itself suggest to make at least 5000[99]. The only models with mutations where pathogencity was highly likely visible were the CYS 62 GLY models that had broken disulfide bridges that could lead to instability in the protein, based on the scores produced by modeller little could be discovered except the higher scores in the disulfide geometry potential (dslf\_fa13) of CYS 62 GLY. Many of the scores show overlap and are difficult to relate with the applied structural changes, it could have been that many of the distributions should overlap more but because few models were made they tend to separate.

The mutated models could have been produced more accurate in the context of making potential disulfide bridges, with the current program they are guessed based on distance which could lead to the insertion of too many or few disulfide bridges. A better method which modeller has to form disulfide bridges based on original data that is available within structures.

HOPE is informative with the results it produces, easy to use, fast and makes structural problems within proteins diagnoseable and understandable when a missense mutation is discovered. However it does not draws a solid conclusion and the information it collects depends on: previous publications, conservation and experimental structures, which will give it a disadvantage when limited knowledge is available. Also it does not asses a complex but it can describe binding sites from the monomer when previously discovered. For the mutation CYS 62 GLY was very clear based on conservation, however for the pathogenic mutation PHE 141 ILE the information was less clear. The provisional likely benign glutamic acid at position 138 makes salt bridges according to HOPE which are strong bonds and can be important for internal structures and binding. With this information GLU 138 ALA would be likely pathogenic, however change in interaction or structure have not been observed with HOPE or SPVAA and therefore can not be classified pathogenic with certainty.

## 6 Conclusion

All used methods use structural information acquired from previous experiments and make it therefore difficult to make predictions based on structural information because many of the structures are fragments.

At its current state VIPUR is not usable for diagnosing defects within proteins and severe changes should be made to make it usable.

SPVAA is only helpful when an expert is available to determine if a mutation has any affect on the protein structure, at its current state it is not user friendly or helpful to inexperienced users, which is reflected in its results. CYS 62 GLY was the only mutation that without much further knowledge could be assumed pathogenic based on the plots and models (Figures 4, 5, 6b, 7b) that were produced, none of the other mutations contained clear information whether they would be pathogenic or not and hardly any differences have been observed between TNF $\alpha$  - $\beta$ .

HOPE is an informative tool that gave new insight in the mutation GLU 138 ALA, in some cases it can be very clear and almost form a conclusion but in other situations it is unable to discover effects of a mutation to elucidate its user and makes the dependence of previously investigated knowledge visible.

## 7 Future work

The VIPUR approach will be further investigated to test whether the features generated by PSI-BLAST are the main predictors. To measure PSI-BLAST feature importance within VIPUR it first needs to be reverse engineered to make it work with Modeller or another tool which is able to implement mutations in PDB files. Once the reverse engineering is finished VIPUR's feature importance is tested with shap[102–108] and/or similar methods and software on the VTS.

SPVAA is highly dependent on available resources and currently comes short to display more potential on the cluster where it ran on, which can be resolved by replacing the Rosetta software which tries to minimize structures of mutated models with other software. Or another possibility is by running SPVAA on a different cluster that has more nodes in its cluster and allows setup jobs that use multiple nodes. SPVAA is not a prediction method and has to be drastically modified to become a variant predictor that is able to predict pathogenicity or deleteriousness, a good starting point for such a predictor would according to the guidelines predicted in section 2.2.

VIPUR and SPVAA could both be improved in various ways, one of them would be by doing molecular dynamic simulations on the mutated structures to determine the effects of structural changes. With SPVAA it would have most likely become clear if the loss of the disulfide bridge from CYS 62 GLY in TNFRSF1A would have caused structural issues. VIPUR could benefit from it as new machine learning feature in situations where limited movement is observed and it changes tremendously when a missense mutation occurred in a protein increases with a mutation or vice versa.

## References

1. Van der Velde, K. J. *et al.* GAVIN: Gene-Aware Variant INterpretation for medical sequencing. *Genome Biology* **18**. ISSN: 1474-7596. doi:10.1186/s13059-016-1141-7. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5240400/> (2019) (Jan. 16, 2017).
2. Baugh, E. H. *et al.* Robust classification of protein variation using structural modelling and large-scale data integration. *Nucleic Acids Research* **44**, 2501–2513. ISSN: 0305-1048 (Apr. 7, 2016).
3. Harper, A. R., Nayee, S. & Topol, E. J. Protective alleles and modifier variants in human health and disease. *Nature Reviews Genetics* **16**, 689–701. ISSN: 1471-0064 (Dec. 2015).
4. NIH. *Sickle cell disease* Genetics Home Reference. <https://ghr.nlm.nih.gov/condition/sickle-cell-disease> (2019).
5. NIH. *Cystic fibrosis* Genetics Home Reference. <https://ghr.nlm.nih.gov/condition/cystic-fibrosis> (2019).
6. NIH, R. G. H. *TRAPS* Genetics Home Reference. <https://ghr.nlm.nih.gov/condition/tumor-necrosis-factor-receptor-associated-periodic-syndrome> (2019).
7. Wikipedia. in *Wikipedia* Page Version ID: 891541555 (Apr. 8, 2019). [https://en.wikipedia.org/w/index.php?title=Protein\\_structure&oldid=891541555](https://en.wikipedia.org/w/index.php?title=Protein_structure&oldid=891541555) (2019).
8. Bennion, B. J. & Daggett, V. Protein Conformation and Diagnostic Tests: The Prion Protein. *Clinical Chemistry* **48**, 2105–2114. ISSN: 0009-9147, 1530-8561 (Dec. 1, 2002).
9. Feyfant, E., Sali, A. & Fiser, A. Modeling mutations in protein structures. *Protein Science : A Publication of the Protein Society* **16**, 2030–2041. ISSN: 0961-8368 (Sept. 2007).
10. Chothia, C. & Lesk, A. M. The relation between the divergence of sequence and structure in proteins. *The EMBO Journal* **5**, 823–826. ISSN: 0261-4189 (Apr. 1986).
11. Ogura, Y. *et al.* A frameshift mutation in NOD2 associated with susceptibility to Crohn's disease. *Nature* **411**, 603. ISSN: 1476-4687 (May 2001).
12. NIH. *The Cost of Sequencing a Human Genome* Genome.gov. <https://www.genome.gov/about-genomics/fact-sheets/Sequencing-Human-Genome-cost> (2019).
13. Hortin, G. L., Carr, S. A. & Anderson, N. L. Introduction: Advances in Protein Analysis for the Clinical Laboratory. *Clinical chemistry* **56**, 149–151. ISSN: 0009-9147 (Feb. 2010).
14. Ng, P. C. & Henikoff, S. SIFT: predicting amino acid changes that affect protein function. *Nucleic Acids Research* **31**, 3812–3814. ISSN: 0305-1048 (July 1, 2003).
15. Nanev, C. N. How do crystal lattice contacts reveal protein crystallization mechanism? *Crystal Research and Technology* **43**, 914–920. ISSN: 1521-4079 (2008).
16. Niu, B. *et al.* Protein-structure-guided discovery of functional mutations across 19 cancer types. *Nature Genetics* **48**, 827–837. ISSN: 1546-1718 (2016).
17. wwPDB. *wwPDB: Deposition Statistics* <http://www.wwpdb.org/stats/deposition> (2019).
18. Cantrill, S. *Chemiotics: How many proteins can we make? : The Sceptical Chymist* [http://blogs.nature.com/thescepticalchymist/2008/04/chemiotics\\_how\\_many\\_proteins\\_c.html](http://blogs.nature.com/thescepticalchymist/2008/04/chemiotics_how_many_proteins_c.html) (2019).
19. PDB101. *PDB101: Learn: Guide to Understanding PDB Data: Missing Coordinates and Biological Assemblies* RCSB: PDB-101. <https://pdb101.rcsb.org/learn/guide-to-understanding-pdb-data/missing-coordinates-and-biological-assemblies> (2019).
20. Ridgen, D. J. in *From Protein Structure to Function With Bioinformatics* 2nd ed., 503 (Springer). ISBN: ISBN 978-94-024-1069-3.
21. Yonath, A. X-ray crystallography at the heart of life science. *Current Opinion in Structural Biology. Carbohydrates and glycoconjugates/Biophysical methods* **21**, 622–626. ISSN: 0959-440X (Oct. 1, 2011).

22. Wikipedia. in *Wikipedia* Page Version ID: 895739876 (May 6, 2019). [https://en.wikipedia.org/w/index.php?title=Ramachandran\\_plot&oldid=895739876](https://en.wikipedia.org/w/index.php?title=Ramachandran_plot&oldid=895739876) (2019).
23. Shourya, S., Burman, R. & Mulligan, V. K. *Scoring Tutorial* <https://rosettacommons.org/demos/latest/tutorials/scoring/scoring#comparing-rosetta-scores-to-real-life-energies> (2019).
24. Andreeva, A., Howorth, D., Chothia, C., Kulesha, E. & Murzin, A. G. SCOP2 prototype: a new approach to protein structure mining. *Nucleic Acids Research* **42**, D310–D314. ISSN: 0305-1048 (Database issue Jan. 1, 2014).
25. Wikipedia. in *Wikipedia* Page Version ID: 898321911 (May 22, 2019). [https://en.wikipedia.org/w/index.php?title=Membrane\\_protein&oldid=898321911](https://en.wikipedia.org/w/index.php?title=Membrane_protein&oldid=898321911) (2019).
26. Wikipedia. in *Wikipedia* Page Version ID: 898320912 (May 22, 2019). [https://en.wikipedia.org/w/index.php?title=Globular\\_protein&oldid=898320912](https://en.wikipedia.org/w/index.php?title=Globular_protein&oldid=898320912) (2019).
27. Wikipedia. in *Wikipedia* Page Version ID: 867647761 (Nov. 7, 2018). <https://en.wikipedia.org/w/index.php?title=Scleroprotein&oldid=867647761> (2019).
28. Wikipedia. in *Wikipedia* Page Version ID: 891043075 (Apr. 5, 2019). [https://en.wikipedia.org/w/index.php?title=Intrinsically\\_disordered\\_proteins&oldid=891043075](https://en.wikipedia.org/w/index.php?title=Intrinsically_disordered_proteins&oldid=891043075) (2019).
29. Stephanie. *Monte Carlo Simulation / Method Statistics How To*. <https://www.statisticshowto.datasciencecentral.com/monte-carlo-simulation/> (2019).
30. Wikipedia. in *Wikipedia* Page Version ID: 896113843 (May 8, 2019). [https://en.wikipedia.org/w/index.php?title=Monte\\_Carlo\\_method&oldid=896113843](https://en.wikipedia.org/w/index.php?title=Monte_Carlo_method&oldid=896113843) (2019).
31. Wikipedia. in *Wikipedia* Page Version ID: 51889441 (July 2, 2018). <https://nl.wikipedia.org/w/index.php?title=Monte-Carlosimulatie&oldid=51889441> (2019).
32. Alon Honig. *Introduction to Monte Carlo Methods* <https://www.youtube.com/watch?v=t0F3S-46bIQ> (2019).
33. Roth-Wojcicki, E. *Tumor Necrosis Factor Receptor Associated Periodic Syndrome (Juvenile)* <https://www.rheumatology.org/I-Am-A/Patient-Caregiver/Diseases-Conditions/Tumor-Necrosis-Factor-Receptor-Associated-Periodic-Syndrome-Juvenile> (2019).
34. Kimberley, F. C., Lobito, A. A., Siegel, R. M. & Screaton, G. R. Falling into TRAPS-receptor misfolding in the TNF receptor 1-associated periodic fever syndrome. *Arthritis Research & Therapy* **9**, 217. ISSN: 1478-6362 (2007).
35. Aksentijevich, I. et al. The tumor-necrosis-factor receptor-associated periodic syndrome: new mutations in TNFRSF1A, ancestral origins, genotype-phenotype studies, and evidence for further genetic heterogeneity of periodic fevers. *American Journal of Human Genetics* **69**, 301–314. ISSN: 0002-9297 (Aug. 2001).
36. Gray, P. W., Barrett, K., Chantry, D., Turner, M. & Feldmann, M. Cloning of human tumor necrosis factor (TNF) receptor cDNA and expression of recombinant soluble TNF-binding protein. *Proceedings of the National Academy of Sciences of the United States of America* **87**, 7380–7384. ISSN: 0027-8424 (Oct. 1990).
37. Walter, R. & Stefan, O. in *Encyclopedia of Molecular Pharmacology* 2nd ed., 1505 (Springer, Nov. 2007). ISBN: 978-3-540-38918-7 978-3-540-38921-7.
38. Banner, D. W. et al. Crystal structure of the soluble human 55 kd TNF receptor-human TNF complex: Implications for TNF receptor activation. *Cell* **73**, 431–445. ISSN: 0092-8674 (May 7, 1993).
39. Segueni, N. et al. Innate myeloid cell TNFR1 mediates first line defence against primary Mycobacterium tuberculosis infection. *Scientific Reports* **6**. ISSN: 2045-2322. doi:10.1038/srep22454. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4773807/> (2019) (Mar. 2, 2016).

40. Naismith, J. H., Devine, T. Q., Brandhuber, B. J. & Sprang, S. R. Crystallographic Evidence for Dimerization of Unliganded Tumor Necrosis Factor Receptor. *Journal of Biological Chemistry* **270**, 13303–13307. ISSN: 0021-9258, 1083-351X (June 2, 1995).
41. Bender, L. M., Morgan, M. J., Thomas, L. R., Liu, Z.-G. & Thorburn, A. The adaptor protein TRADD activates distinct mechanisms of apoptosis from the nucleus and the cytoplasm. *Cell Death & Differentiation* **12**, 473. ISSN: 1476-5403 (May 2005).
42. Muppidi, J. R., Tschopp, J. & Siegel, R. M. Life And Death Decisions: Secondary Complexes and Lipid Rafts in TNF Receptor Family Signal Transduction. *Immunity* **21**, 461–465. ISSN: 1074-7613 (Oct. 1, 2004).
43. Vinay, K., Abul, K. & Jon, C. in *Robbins and Cotran Pathologic Basis of Disease, Professional Edition* 9th, 1464 (Elsevier, July 9, 2014). ISBN: 978-0-8153-4432-2.
44. Chen, G. & Goeddel, D. V. TNF-R1 Signaling: A Beautiful Pathway. *Science* **296**, 1634–1635. ISSN: 0036-8075, 1095-9203 (May 31, 2002).
45. Hengartner, M. O. The biochemistry of apoptosis. *Nature* **407**, 770. ISSN: 1476-4687 (Oct. 2000).
46. Aggarwal, B. B., Eessalu, T. E. & Hass, P. E. Characterization of receptors for human tumour necrosis factor and their regulation by -interferon. *Nature* **318**, 665. ISSN: 1476-4687 (Dec. 1985).
47. Hamosh, A. & McKusick, V. A. OMIM Entry - \* 153440 - LYMPHOTOXIN-ALPHA; LTA <https://omim.org/entry/153440?search=lymphotoxin&highlight=lymphotoxin> (2019).
48. Kriegler, M., Perez, C., DeFay, K., Albert, I. & Lu, S. D. A novel form of TNF/cachectin is a cell surface cytotoxic transmembrane protein: Ramifications for the complex physiology of TNF. *Cell* **53**, 45–53. ISSN: 0092-8674 (Apr. 8, 1988).
49. Pieper, U. *et al.* modbase, a database of annotated comparative protein structure models and associated resources. *Nucleic Acids Research* **37**, D347–D354. ISSN: 0305-1048 (Database issue Jan. 2009).
50. Guex, N., Peitsch, M. C. & Schwede, T. Automated comparative protein structure modeling with SWISS-MODEL and Swiss-PdbViewer: A historical perspective. *ELECTROPHORESIS* **30**, S162–S173. ISSN: 1522-2683 (S1 2009).
51. Bertoni, M., Kiefer, F., Biasini, M., Bordoli, L. & Schwede, T. Modeling protein quaternary structure of homo- and hetero-oligomers beyond binary interactions by homology. *Scientific Reports* **7**, 10480. ISSN: 2045-2322 (Sept. 5, 2017).
52. Benkert, P., Biasini, M. & Schwede, T. Toward the estimation of the absolute quality of individual protein structure models. *Bioinformatics* **27**, 343–350. ISSN: 1367-4803 (Feb. 1, 2011).
53. Bienert, S. *et al.* The SWISS-MODEL Repository—new features and functionality. *Nucleic Acids Research* **45**, D313–D319. ISSN: 0305-1048 (D1 Jan. 4, 2017).
54. Waterhouse, A. *et al.* SWISS-MODEL: homology modelling of protein structures and complexes. *Nucleic Acids Research* **46**, W296–W303. ISSN: 0305-1048 (W1 July 2, 2018).
55. Poultney, C. S. *et al.* Rational Design of Temperature-Sensitive Alleles Using Computational Structure Prediction. *PLOS ONE* **6**, e23947. ISSN: 1932-6203 (Sept. 2, 2011).
56. Baugh, E. H. *et al.* SUPPLEMENTARY: Robust classification of protein variation using structural modelling and large-scale data integration. *Nucleic Acids Research* **44**, 2501–2513. ISSN: 0305-1048 (Apr. 7, 2016).
57. Commons, R. *About — RosettaCommons* <https://www.rosettacommons.org/about> (2019).
58. Commons, R. *How to prepare structures for use in Rosetta* [https://www.rosettacommons.org/docs/latest/rosetta\\_basics/preparation/preparing\\_structures](https://www.rosettacommons.org/docs/latest/rosetta_basics/preparation/preparing_structures) (2019).
59. Commons, R. *Relax application* [https://www.rosettacommons.org/docs/latest/application\\_documentation/structure\\_prediction/relax](https://www.rosettacommons.org/docs/latest/application_documentation/structure_prediction/relax) (2019).

60. Conway, P., Tyka, M. D., DiMaio, F., Konerding, D. E. & Baker, D. Relaxation of backbone bond geometry improves protein energy landscape modeling. *Protein Science : A Publication of the Protein Society* **23**, 47–55. ISSN: 0961-8368 (Jan. 2014).
61. Tyka, M. D. *et al.* Alternate States of Proteins Revealed by Detailed Energy Landscape Mapping. *Journal of Molecular Biology* **405**, 607–618. ISSN: 0022-2836 (Jan. 14, 2011).
62. Leaver-Fay, A. & Kellogg, E. *ddg\_monomer application* [https://www.rosettacommons.org/docs/latest/application\\_documentation/analysis/ddg-monomer](https://www.rosettacommons.org/docs/latest/application_documentation/analysis/ddg-monomer) (2019).
63. Jared, A.-B. *Score Commands* [https://www.rosettacommons.org/docs/latest/application\\_documentation/analysis/score-commands](https://www.rosettacommons.org/docs/latest/application_documentation/analysis/score-commands) (2019).
64. Betancourt, M. R. Efficient Monte Carlo trial moves for polypeptide simulations. *The Journal of Chemical Physics* **123**, 174905. ISSN: 0021-9606 (Oct. 31, 2005).
65. Smith, C. A. *Backrub application* [https://www.rosettacommons.org/docs/latest/application\\_documentation/structure\\_prediction/backrub](https://www.rosettacommons.org/docs/latest/application_documentation/structure_prediction/backrub) (2019).
66. Jeffrey, J. G., Sergey, L. & Team, P. *PyRosetta* <http://www.pyrosetta.org/> (2019).
67. NCBI. *PSIBLAST* <http://www.biology.wustl.edu/gcg/psiblast.html> (2019).
68. NCBI. *PSSM Viewer* [https://www.ncbi.nlm.nih.gov/Class/Structure/pssm/pssm\\_viewer.cgi](https://www.ncbi.nlm.nih.gov/Class/Structure/pssm/pssm_viewer.cgi) (2019).
69. Wikipedia. in *Wikipedia* Page Version ID: 890694148 (Apr. 2, 2019). <https://en.wikipedia.org/w/index.php?title=BLAST&oldid=890694148> (2019).
70. Word, J. M. *et al.* Visualizing and quantifying molecular goodness-of-fit: small-probe contact dots with explicit hydrogen atoms<sup>11</sup>Edited by J. Thornton. *Journal of Molecular Biology* **285**, 1711–1733. ISSN: 0022-2836 (Jan. 29, 1999).
71. LAB, R. *Probe Software : Kinemage Website* <http://kinemage.biochem.duke.edu/software/probe.php> (2019).
72. Song, Y. *et al.* High-Resolution Comparative Modeling with RosettaCM. *Structure* **21**, 1735–1742. ISSN: 0969-2126 (Oct. 8, 2013).
73. Soding, J. Protein homology detection by HMM-HMM comparison. *Bioinformatics* **21**, 951–960. ISSN: 1367-4803, 1460-2059 (Apr. 1, 2005).
74. Källberg, M. *et al.* Template-based protein structure modeling using the RaptorX web server. *Nature protocols* **7**, 1511–1522. ISSN: 1754-2189 (July 19, 2012).
75. Yang, Y., Faraggi, E., Zhao, H. & Zhou, Y. Improving protein fold recognition and template-based modeling by employing probabilistic-based matching between predicted one-dimensional structural properties of query and corresponding native properties of templates. *Bioinformatics* **27**, 2076–2082. ISSN: 1367-4803 (Aug. 1, 2011).
76. Ovchinnikov, S. *et al.* Protein structure determination using metagenome sequence data. *Science* **355**, 294–298. ISSN: 0036-8075, 1095-9203 (Jan. 20, 2017).
77. LAB, Z. *LOMETS* <https://zhanglab.ccmb.med.umich.edu/LOMETS/help.html> (2019).
78. Wu, S. & Zhang, Y. LOMETS: A local meta-threading-server for protein structure prediction. *Nucleic Acids Research* **35**, 3375–3382. ISSN: 0305-1048 (May 2007).
79. Modeller. *About MODELLER* <https://salilab.org/modeller/> (2019).
80. Eswar, N. *et al.* Comparative Protein Structure Modeling Using Modeller. *Current protocols in bioinformatics / editorial board, Andreas D. Baxevanis ... [et al.]* **0 5**, Unit-5.6. ISSN: 1934-3396 (Oct. 2006).
81. ensembl. *Variant Effect Predictor - Homo sapiens - GRCh37 Archive browser 96* [http://grch37.ensembl.org/Homo\\_sapiens/Tools/VEP](http://grch37.ensembl.org/Homo_sapiens/Tools/VEP) (2019).

82. NCBI. *Representation of clinical significance in ClinVar and other variation resources at NCBI* <https://www.ncbi.nlm.nih.gov/clinvar/docs/clinsig/> (2019).
83. gnomAD. *gnomAD* <https://gnomad.broadinstitute.org/> (2019).
84. Sarrauste de Menthière, C. *et al.* INFEVERS: the Registry for FMF and hereditary inflammatory disorders mutations. *Nucleic Acids Research* **31**, 282–285. ISSN: 0305-1048 (Jan. 1, 2003).
85. Naismith, J. H., Devine, T. Q., Kohno, T. & Sprang, S. R. Structures of the extracellular domain of the type I tumor necrosis factor receptor. *Structure* **4**, 1251–1262. ISSN: 0969-2126 (Nov. 15, 1996).
86. Burley, S. K. *et al.* RCSB Protein Data Bank: Sustaining a living digital data resource that enables breakthroughs in scientific research and biomedical education. *Protein Science* **27**, 316–330. ISSN: 0961-8368 (Jan. 1, 2018).
87. Consortium, T. U. UniProt: a hub for protein information. *Nucleic Acids Research* **43**, D204–D212. ISSN: 0305-1048 (Database issue Jan. 28, 2015).
88. Schrödinger. *PyMOL — pymol.org* <https://pymol.org/2/> (2019).
89. Venselaar, H., te Beek, T. A., Kuipers, R. K., Hekkelman, M. L. & Vriend, G. Protein structure analysis of mutations causing inheritable diseases. An e-Science approach with life scientist friendly interfaces. *BMC Bioinformatics* **11**, 548. ISSN: 1471-2105 (Dec. 2010).
90. CMBI. *HOPE* <http://www.cmbi.ru.nl/hope/> (2019).
91. CMBI. *HOPE about* <http://www.cmbi.ru.nl/hope/about/> (2019).
92. CMBI. *HOPE methods* <http://www.cmbi.ru.nl/hope/method/> (2019).
93. Wickham, H. *Create Elegant Data Visualisations Using the Grammar of Graphics* <https://ggplot2.tidyverse.org/> (2019).
94. Dowle, m., Srinivasan, A., Gorecki, J. & Chirico, M. *R’s data.table package extends data.frame: Contribute to Rdatatable/data.table development by creating an account on GitHub* original-date: 2014-06-07T16:38:05Z. June 11, 2019. <https://github.com/Rdatatable/data.table> (2019).
95. Baugh, E. H. VIPUR: Variant Interpretation and Prediction Using Rosetta. doi:None. <https://osf.io/bd2h4/> (2019) (Sept. 15, 2015).
96. Sukits, S. F. *et al.* Solution structure of the tumor necrosis factor receptor-1 death domain. *Journal of Molecular Biology* **310**, 895–906. ISSN: 0022-2836 (July 20, 2001).
97. Liu, G.-H. *et al.* Lipin proteins form homo- and hetero-oligomers. *The Biochemical journal* **432**, 65–76. ISSN: 0264-6021 (Oct. 25, 2010).
98. Koshland, D. E. Application of a Theory of Enzyme Specificity to Protein Synthesis\*. *Proceedings of the National Academy of Sciences of the United States of America* **44**, 98–104. ISSN: 0027-8424 (Feb. 1958).
99. Commons, R. *Analyzing Results* <https://www.rosettacommons.org/docs/latest/getting-started/Analyzing-Results> (2019).
100. Li, S. C., Goto, N. K., Williams, K. A. & Deber, C. M. Alpha-helical, but not beta-sheet, propensity of proline is determined by peptide environment. *Proceedings of the National Academy of Sciences of the United States of America* **93**, 6676–6681. ISSN: 0027-8424 (June 25, 1996).
101. Ambrogelly, A., Paloura, S. & Söll, D. Natural expansion of the genetic code. *Nature Chemical Biology* **3**, 29–35. ISSN: 1552-4469 (Jan. 2007).
102. Slundberg. *slundberg/shap: A unified approach to explain the output of any machine learning model.* <https://github.com/slundberg/shap> (2019).
103. Štrumbelj, E. & Kononenko, I. Explaining Prediction Models and Individual Predictions with Feature Contributions. *Knowl. Inf. Syst.* **41**, 647–665. ISSN: 0219-1377 (Dec. 2014).

104. Ribeiro, M. T., Singh, S. & Guestrin, C. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. *arXiv:1602.04938 [cs, stat]*. arXiv: 1602.04938. <http://arxiv.org/abs/1602.04938> (2019) (Feb. 16, 2016).
105. Shrikumar, A., Greenside, P. & Kundaje, A. *Learning Important Features Through Propagating Activation Differences* in International Conference on Machine Learning International Conference on Machine Learning (July 17, 2017), 3145–3153. <http://proceedings.mlr.press/v70/shrikumar17a.html> (2019).
106. Datta, A., Sen, S. & Zick, Y. Algorithmic Transparency via Quantitative Input Inuence: 20.
107. Bach, S. *et al.* On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Relevance Propagation. *PLoS ONE* **10**. ISSN: 1932-6203. doi:10.1371/journal.pone.0130140. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4498753/> (2019) (July 10, 2015).
108. Datadive. *Interpreting random forests — Diving into data* Datadive. <http://blog.datadive.net/interpreting-random-forests/> (2019).

## **Supplementary**