

# Attentive Neural Protein Structure Prediction

---

Prediction of protein tertiary structure via the translation of amino acid sequences into internal torsional angles.

College of Engineering, Mathematics and Physical Sciences  
The University of Exeter

1 May 2019

## Abstract

Proteins are abundant in virtually all biological processes within the cells of, not just humans, but all living organisms. Sequences of amino acids undergo a spontaneous transition into their final 3-dimensional structure, called folding. Structure determination is important as the structure dictates the function of the protein. This report devises a novel method of protein structure prediction that draws on recent successes in the fields of Neural Machine Translation and Natural Language Processing, applying them to the field of Bioinformatics. The proposed method is a deep attentive encoder-decoder network that translates embedded ngrams of amino acid sequences into torsion angles, feeding them into the pNeRF algorithm for a Cartesian representation. The results are then benchmarked against the contestants of CASP13.

*I certify that all material in this dissertation which is not my own work has been identified.*

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Background</b>	<b>3</b>
2.1	The Folding Process . . . . .	3
2.2	Bioinformatics Background . . . . .	5
<b>3</b>	<b>Related Research</b>	<b>6</b>
3.1	Anton . . . . .	6
3.2	I-Tasser . . . . .	6
3.3	AlphaFold . . . . .	7
3.4	Google Neural Machine Translation . . . . .	7
<b>4</b>	<b>Model</b>	<b>7</b>
4.1	Dataset . . . . .	7
4.2	Representations . . . . .	8
4.2.1	Input Representation . . . . .	8
4.2.2	Output Representation . . . . .	10
4.3	Data Pipeline . . . . .	11
4.4	The Network . . . . .	12
4.4.1	Optimisation . . . . .	12
4.4.2	Residual Connections . . . . .	15
4.4.3	Gradient Clipping . . . . .	16
4.4.4	Dropout . . . . .	17
4.5	Attention . . . . .	18
<b>5</b>	<b>Implementation</b>	<b>19</b>
5.1	Inference . . . . .	20
5.2	Torsion Space to Cartesian Space . . . . .	20
5.3	Global Distance Test Total Score . . . . .	20
<b>6</b>	<b>Results</b>	<b>21</b>
<b>7</b>	<b>Analysis</b>	<b>22</b>
7.1	Angle Distributions . . . . .	22
7.2	Reconstruction . . . . .	23
7.3	Attentiveness . . . . .	24
7.4	A Landscape Comparison . . . . .	25
<b>8</b>	<b>Future Work</b>	<b>26</b>
<b>9</b>	<b>Conclusion</b>	<b>27</b>
<b>A</b>	<b>PDB Codes</b>	<b>33</b>

<b>B</b>	<b>Embeddings</b>	<b>34</b>
B.1	1-Grams . . . . .	34
B.2	2-Grams . . . . .	35
B.3	3-Grams . . . . .	36

# 1 Introduction

Proteins are ubiquitous in a vast number of biological processes in the cells of all organic organisms. They have various functions ranging from antibodies that bind to foreign particles, such as viruses and bacteria, enzymes that act as catalysts in chemical reactions within cells or messengers that transmit signals to coordinate biological processes [1]. They also provide structural support for cells and can transport and store atoms and molecules within them. Computational models of protein folding are vital as the structure of a folded protein dictates its function. Protein structure prediction is an active field of research that remains elusive, captivating the minds of many researchers for over half a century. The fundamental principles of folding have practical applications in the exploitation of the advances in genome research, in the understanding of different pathologies and in the design of novel proteins with special functions, just to name a few.

The costly nature of experimental methods of protein structure determination forms the necessity to develop accurate and efficient computational models instead. Since its inception in 1994, Critical Assessment of protein Structure Prediction (CASP) has run every 2 years for precisely this reason. The competition is conducted in a double-blind fashion; neither the predictors nor the organizers and assessors know the structures of the target protein at the time the predictions are made. This provides the benchmark for comparing the predictive methods and is how the model devised in this report will be assessed.

Currently, no explicit sequence-to-structure mapping exists that is scalable. Molecular Dynamics simulations are accurate but inefficient and Template Modelling is efficient but inaccurate where no close homolog exists. However, Google DeepMind entered CASP for the first time this year with AlphaFold. Much to the surprise of the regular contestants, this deep neural network approach trounced the competition by a statistically significant margin. The model devised in this report intends to capitalise on the success of deep neural networks applied to protein structure prediction, taking inspiration from AlphaFold and other deep learning approaches. The objectives of this research are to conduct a *de novo* predictor that is faster than Molecular Dynamics simulations and more accurate than *ab initio* Template Based Modelling.

The report frames the task of protein structure prediction as a Neural Machine Translation Problem, taking an embedded representation of a sequence of amino acids that have been split into ngrams and translating them into a sequence of discretized internal torsional angles, with an attentive encoder-decoder. The predicted sequence is then fed into the parallelised Natural extension Reference Frame (pNeRF) algorithm to convert the dihedral angles into Cartesian coordinates. The paper shows the promise of sequence-to-sequence models for protein folding, benchmarking the results against those of the contenders of CASP13.

# 2 Background

## 2.1 The Folding Process

Under physiological conditions, a protein undergoes a spontaneous transition called folding. Proteins are linear polymers built from a chain of 20 possible amino acids that consist of backbone atoms and a side chain, which differentiates them. The side chain determines the hydrophobicity, charge and polarity of the amino acid, which are all vital physiochemical features in the folding process. There are three distinct stages to the folding process, identified by the structures they create.

Primary structure is the first stage. Chains of amino acids are linked by peptide bonds, referred to as polypeptides. The polypeptide backbone is a repeating sequence along the polypeptide chain. Peptide bonds form through the nucleophilic addition-elimination reaction between the carboxyl group of one amino acid and the amino group of another. X-ray diffraction has found that these peptide bonds are rigid and planar (Figure 2). Thus, for a pair of adjacent amino acids, six atoms lie in the same plane. As discussed later, this provides a useful simplification to the tertiary structure representation. The start of a polypeptide will always be a nitrogen atom and the end of a polypeptide will be a carbon atom; the N-terminus and C-terminus, respectively. Once in the chain, each amino acid is referred to as a residue.

The next step in the folding process is typically towards two motifs,  $\alpha$ -helices and  $\beta$ -sheets. Secondary structure is primarily determined by backbone interactions such as hydrogen bonding.  $\alpha$ -helices are formed by the hydrogen bonding of the backbone into spiral shapes where the hydrogen bonds run up and down, stabilising the structure. There are two kinds of  $\beta$ -sheets. When the N-terminus and C-terminus are aligned in the same direction it is a parallel  $\beta$ -sheet and when they alternate it is an anti-parallel  $\beta$ -sheet, or  $\beta$ -turn.

Finally, the tertiary structure of a protein is also stabilised by hydrogen bonds but new interactions are now introduced to the folding process. The major driving force is the polarity of a residues side chain. Polarity determines how hydrophobic or hydrophilic the residue is. The hydrophobic effect dictates that in aqueous environments, such as cells, non-polar side chains are hydrophobic meaning they tend to cluster in the centre of globular proteins, away from the aqueous surroundings. On the other hand, polar side chains favour interactions with the surrounding solution and so tend towards the surface of a protein [2]. There are five polar amino acids that can carry a charge indicated by the residues pKa value that are split into two groups, basic and acidic. Basic bear a full positive charge while acidic bear a full negative charge, at the normal physiological pH. If an acidic residue comes into close proximity with a basic residue during the folding process, ionic bonds form as they have opposing charges. The van der Waals force is an attractive force that balances the repulsion that increases when two non-polar atoms approach each other. As a result, there is a distance at which repulsive and attractive forces precisely balance, stabilising the hydrophobic core of proteins [3]. Disulphide bonds are sulphur-sulphur covalent bonds that form between the thiol groups of two cysteine residues. These are very strong intra-molecular bonds that play an important role in determining the structure of a protein. Section 4.2 displays how the physiochemical properties are implicitly encoded in the input representation via an embedding in vector space.

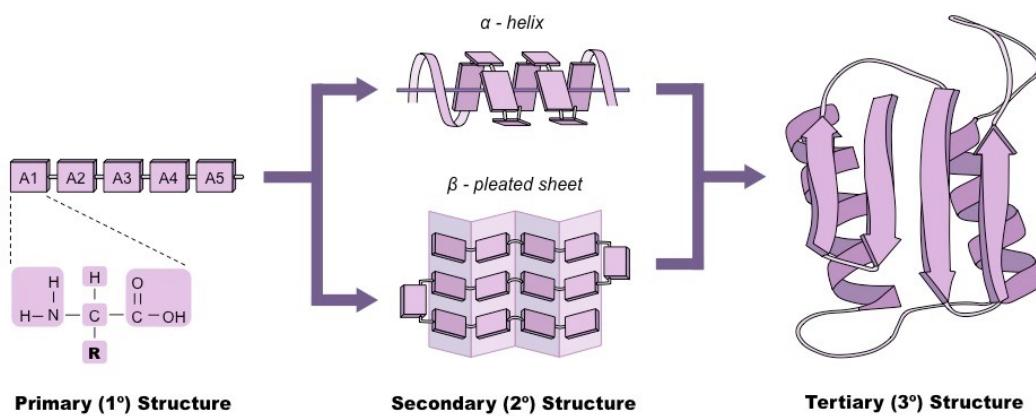


Figure 1: The Folding Process

## 2.2 Bioinformatics Background

Before 1963, various types of polypeptide chain configurations had been proposed, notably  $\alpha$ -helices. However, there was no analytical method of writing these configurations until the physicist, G. N. Ramachandran, devised the Ramachandran Plot (Figure 4, Figure 13) [4]. The peptide bond between the carbonyl  $C$  and  $N$  only exists in two discrete configurations, trans and cis isomers. 99.9% of peptide bonds adopt the trans isomer with cis isomers reflecting trans isomers at exactly  $\pi$  radians, denoted by the  $\omega$  angle. The  $\phi$  and  $\psi$  angles are the angle of rotation about the single bonds  $N - C_\alpha$  and  $C_\alpha - C$  respectively. Plotting  $\phi$  and  $\psi$  for all residues in a chain show the allowed confirmations of proteins whilst also indicating secondary structure [5]. It is clear from the Ramachandran plots that not all rotations are allowed, due to steric hindrance.

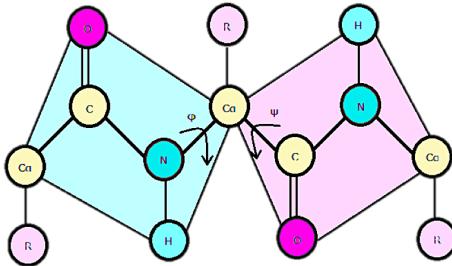


Figure 2: Peptide Bond

In 1973, Anfinsen conducted a Nobel prize-winning experiment that lead to what is now referred to as the Thermodynamic Hypothesis. The hypothesis states that the three-dimensional structure of a native protein in its normal physiological milieu (solvent, pH, ionic strength, presence of other components such as metal ions or prosthetic groups, temperature and other) is the one in which the Gibbs free energy of the whole system is lowest; that is, that the native confirmation is determined by the totality of inter-atomic interactions and hence by the amino acid sequence, in a given environment [6]. Throughout the 1950s, Anfinsen performed a series of unscrambling experiments in support of his eventual hypothesis, that showed primary structure determines the confirmation of a protein [7, 8]. Denaturing agents were mixed with Ribonuclease to create the denatured enzyme by breaking its 4 disulphide bonds. When removing one of the denaturing agents the enzyme refolded to an inactive scrambled enzyme with incorrect disulphide bonds. Upon adding back trace amounts of the denaturing agent, the incorrect disulphide bonds broke causing it to eventually refold to its native state. This showed that for small globular proteins in their standard physiological environment, native structure is determined purely by the proteins amino acid sequence [9].

Several experimental methods are currently used to determine protein structure, including X-ray crystallography and NMR spectroscopy. In most cases, experimental data is not sufficient and bond lengths and angles are often added manually resulting in structures comprised of a balanced mixture of experimental observation and knowledge-based modelling. As these experimental structures are used as the training data it is important to be aware of this.

Upon determining protein structure from X-ray crystallography, the protein is purified, crystallised and subjected to an intense beam of X-rays [10]. The crystal diffracts the beam into a pattern, showing the distribution of electrons in the protein. The electron density map can be interpreted to provide information about the location of each atom in the chain. Although X-ray crystallography shows detailed information about the protein, crystallisation can be difficult.

Flexible proteins do not form useful electron density maps like rigid proteins because crystallography requires many molecules to be aligned in the same orientation [11]. The X-ray resolution is a metric that measures the accuracy of crystallographic structures. The training data for the proposed model consists of proteins whose structure was determined via X-ray crystallography. This is discussed further in section 4.

These methods are expensive and time consuming which means great care must be taken when allocating resources. Given this expense, determining protein structure *in silico* is preferable.

## 3 Related Research

### 3.1 Anton

Molecular Dynamics (MD) simulations can be used to model molecular systems such as proteins with an atomic level of precision. Many biological processes occur on millisecond timescales yet MD simulations on this scale lie beyond the reach of current general-purpose technology, by several orders of magnitude. D.E. Shaw Research is a biochemical research company engaged in the design of novel algorithms and machine learning architectures for high-speed MD simulations of proteins and other macromolecules.

Anton is their specialised, massively parallel machine designed to accelerate MD simulations [12]. It simulates atomic motion over a period of time according to classical physics. At each discrete time step the force on each particle due to other particles is computed and the net force is used to update each particles position and velocity. The forces computed are bond, van der Waals and electrostatic forces. The decrease in simulation time comes from specialised hardware that outperforms the expected speedup predicted by Moores Law.

Although Anton has enabled the longest MD simulations to date by a very large factor, its performance is dependent on the size of the system being simulated. For instance, the protein BPTI took  $18.2 \mu\text{s}/\text{day}$  to simulate, which has over 17,000 particles [13]. Despite these speed ups, Anton does not scale to long sequences or slow folding proteins. The Folding@home project found NTL9, the slowest-folding protein folded *ab initio* by MD simulation. NTL9 has a folding time of 1.5ms which means Anton would take approximately 3 months to simulate this singular protein [14].

### 3.2 I-Tasser

A typical template-based modelling procedure involves two major steps: finding proteins with sequences similar to known structure(s), and building 3D models using the detected homologues as structural templates. There were notable improvements from CASP11 to CASP12 in the template-based modelling category with two servers from the Zhang group, Zhang-server and QUARK, outperforming the rest of the servers in a statistically significant manner [15]. They placed second in CASP13 when looking at overall results.

Prior to CASP13, I-TASSER (Zhang-server) has consistently ranked first in CASP competitions [16]. Its name comes from its hierarchical approach of: Threading, ASSEmbly and Refinement. Threading refers to the bioinformatics procedure of identifying template proteins from structure databases. PSI-BLAST matches a query sequence to evolutionary relatives from the database. A sequence profile is created from Multiple-Sequence Alignments (MSA) and used to predict secondary structure using PSIPRED. This is then fed through seven state-of-the-art threading

programs with the top templates from each selected for assembly. Structural assembly involves continuous fragments being excised from the templates and used to build confirmations of well-aligned segments, with unaligned regions being built *ab initio*. The fragment assembly is performed by a Monte-Carlo simulation, with this process being iteratively performed for refinement.

Template-based methods suffer when sequence similarity falls below 20% as the *ab initio* modelling is highly inaccurate.

### 3.3 AlphaFold

Googles DeepMind entered CASP13 with AlphaFold, their first attempt at protein structure prediction that achieved surprisingly unprecedeted results [17]. DeepMind placed themselves a comfortable distance from the second place predictor (the Zhang group) in the free modelling category, which focusses on modelling novel protein folds.

AlphaFold uses MSAs and profiles generated from HHBlits [18] and PSI-BLAST as its input, and a SoftMax over discretized spatial ranges as its output, predicting a probability distribution over inter-residue distances. The distribution from a deep residual convolutional neural network trained on a non-redundant database of proteins selected from the PDB is used as a statistical potential function that is minimised to generate the protein fold. The potential is normalised using a learned reference state and combined with a more traditional physics-based potential. Minimisation of the combined energy function uses L-BFGS gradient descent to yield high accuracy folds.

As the paper is yet to be released, at the time of writing, a more comprehensive analysis of AlphaFold is not possible.

### 3.4 Google Neural Machine Translation

Away now from computational protein folding models to Google Neural Machine Translation (GNMT), the algorithm that inspired framing the task of protein structure prediction as a NMT problem [19]. In 2016 Google switched their translate service over to the GNMT program to overcome the weaknesses of their conventional phrase-based translation systems.

The new model consists of a deep LSTM network with 8 encoder layers and 8 decoder layers using attention and residual connections. The bottom layer of the encoder network is a bidirectional LSTM, with subsequent layers being unidirectional. The decoder layers are all unidirectional. The attention mechanism connects the bottom layer of the decoder to the top layer of the encoder. Words are divided into sub-words, called wordpieces that help remedy problems attributed to the translation of rare words by providing a balance between character-delimited models and word-delimited models. When decoding the output during inference, a beam search technique is utilised that employed length-normalisation.

Using a human-rated side-by-side comparison as a metric, they show that their GNMT system approaches the accuracy achieved by average bilingual human translators.

## 4 Model

### 4.1 Dataset

The Research Collaboratory for Structural Bioinformatics (RCSB) Protein Data Bank (PDB) is a provider of 3D structure data for large biological molecules (proteins, DNA and RNA) typically

obtained via X-ray crystallography, NMR spectroscopy, or, increasingly, cryo-electron microscopy [20]. The search interface enables complex queries of combinations of specific categories. The specific query for this dataset is: sequence identity less than 30%, X-ray resolution less than 2.0Å, chain length between 200 and 400 residues, macromolecular type being protein with no DNA/RNA, singular chains in both the biological assembly and assymmetric unit and a singular entity. The query returns a list of PDB IDs (Appendix A), which are 4-character unique identifiers of entries in the PDB. The BioPython python library uses these IDs to retrieve all the necessary sequence and structure information.

As mentioned previously an objective of this paper is to devise a novel approach to *de novo* protein structure prediction. For this reason, sequence similarity must be low so as to not rely on homologous proteins from the training data when making predictions, constructing a more explicit sequence-to-structure map instead. X-ray resolution below 2.0Å is specified to ensure only well-resolved proteins are included as this implies that the backbone and side chains are clear and hence a more exact protein structure will be learnt by the model. The chain length is above 200 to ensure that clear secondary structure appears in each protein and below 400 to reduce the amount of padding per protein. As the scope of this research is on protein structure prediction the macromolecular type is limited to protein only. Likewise, only singular chains are used as this limits the task solely to the prediction of structure and avoids complexities that arise from ligand binding or docking.

Comprehensive statistical analyses of amino acid sequences in proteins have been conducted by many biologists over the years. As seen in the following section, to qualitatively assess the input embeddings, various physiochemcial residue indexes were used. As a measure of polarity Zimmerman *et al.* took a rough approximation of the electric force due to the different side chains acting on their surroundings [21]. The hydrophobicity index for each residue was obtained by Tanford by collecting solubility data [22]. The normalised van der Waals volume index for each residue was obtained by Fauchere [23]. Likewise, normalised frequencies of secondary structural motifs were obtained by Chou and Fasman from their computerised Chou-Fasman predictive method [24].

## 4.2 Representations

### 4.2.1 Input Representation

Natural Language Processing systems traditionally treat words as discrete atomic symbols. The encodings are arbitrary, providing no useful information to the system regarding any relationships that may exist between the individual symbols. Representing words in this way leads to data sparsity and requires more data in order to successfully train statistical models. To combat this, vector space models embed words in a continuous vector space where semantically similar words appear close to one and other. Embeddings depend on the Distributional Hypothesis, which states “linguistic items with similar distributions have similar meaning” [25].

The word2vec Skip-gram model is a computationally-efficient predictive model for learning word vector representations that are good at predicting nearby words [26]. More formally, given a sequence of training words  $w_1, w_2, w_3, \dots, w_T$ , the objective of the Skip-gram model is to maximise the average log probability:

$$\text{argmax} \frac{1}{T} \sum_{t=0}^T \sum_{\substack{-c < j < c \\ j \neq 0}} \log P(w_{t+j} | w_t) \quad (1)$$

Where  $2c$  is the size of the context window with  $w_t$  as the center,  $T$  is the sequence length. The basic Skip-gram formulation defines  $P(w_{t+j}|w_t)$  using the SoftMax function:

$$P(w_O|w_I) = \frac{\exp(v'^\top_{w_O} v_{w_I})}{\sum_{w=1}^W \exp(v'^\top_w v_{w_I})} \quad (2)$$

Where  $v_w$  and  $v'_w$  are the input and output vector representations of  $w$  and  $W$  is the vocabulary size. This formulation is impractical as the cost of computing  $\nabla \log p(w_O|w_I)$  is proportional to  $W$ . For this reason, Negative Sampling was used instead. Negative Sampling distinguishes the target word  $w_O$  from draws from a noise distribution  $P_n(w)$  using logistic regression, where there are  $k$  negative samples for each data sample and  $\sigma(x) = 1/(1 + \exp(x))$ . Each  $\log P(w_O|w_I)$  is replaced with the following.

$$\log \sigma(v'^\top_{w_O} v_{w_I}) + \sum_{i=1}^k \mathbf{E}_{w_i \sim P_n(w)} [\log(-v'^\top_{w_O} v_{w_I})] \quad (3)$$

This report takes inspiration from ProtVec, which conceptualises that nature uses certain molecular languages to describe biological sequences, such as proteins [27]. Asgari and Mofrad proposed the use of a dense distributed representation of biological sequences to discover the functions encoded within them. In order to train these distributed representations, the amino acid sequences must be broken into subsequences of length  $n$  (i.e. biological words). Typically, n-gram modelling in protein informatics uses an overlapping window of residues. However, ProtVec generates  $n$  lists of shifted non-overlapping biological words instead, resulting in a corpus consisting of  $n \times T$  sequences. The vector representation of the biological words in these sequences are learnt through the Skip-gram neural network. To train the embedding vectors a vector size of 128 is considered and a context window of 20 for values of  $n = 1, 2, 3$ . Thus, each n-gram is presented as a 128-dimensional vector.

To qualitatively analyse the distribution of n-grams within the vector space, TensorFlows TensorBoard Projector is utilised. All n-grams are projected from 128-dimensional space to a 2-dimensional space using t-Distributed Stochastic Neighbour Embedding (t-SNE) [28]. The physiochemical properties of hydrophobicity, polarity and van der Waals volume were analysed alongside the structural motifs of  $\alpha$ -helix,  $\beta$ -sheet and  $\beta$ -turn; using the physiochemical properties from Section 4.1. As there are 20 amino acids, it follows that there is a maximum of 20 1-grams, 400 2-grams and 8000 3-grams in each respective corpus vocabulary (not including start, end or pad tokens). For each n-gram in each vocabulary, the physiochemical properties were averaged based on the individual residues that make up the n-gram and TensorBoard colours them according to this averaged value.

For the small number of 1-grams, TensorBoard is limited to unlabelled categorical colouring, which is of limited value (Appendix B.1). The opposite holds for 3-grams, where the larger vocabulary size meant meaningful clusters failed to form, apart from 3-grams containing cysteine separating themselves from the rest of the vocabulary (Appendix B.3). Figure 3 displays that physiochemical and structural information is well-encoded implicitly for the embedded 2-grams. Polar uncharged 2-grams cluster together (lower-right), typically away from hydrophobic 2-grams (upper-left). Likewise, larger van der Waals volumes favour the upper-right of the spherical t-SNE embedding. As the hydrophobicity, polarity and van der Waals volume affect the secondary structure formation, observing the conservation of clusters is promising. Consequently, structure is implicitly encoded aswell.  $\beta$ -sheets favour the opposite space to  $\beta$ -turns and  $\alpha$ -helices encompass

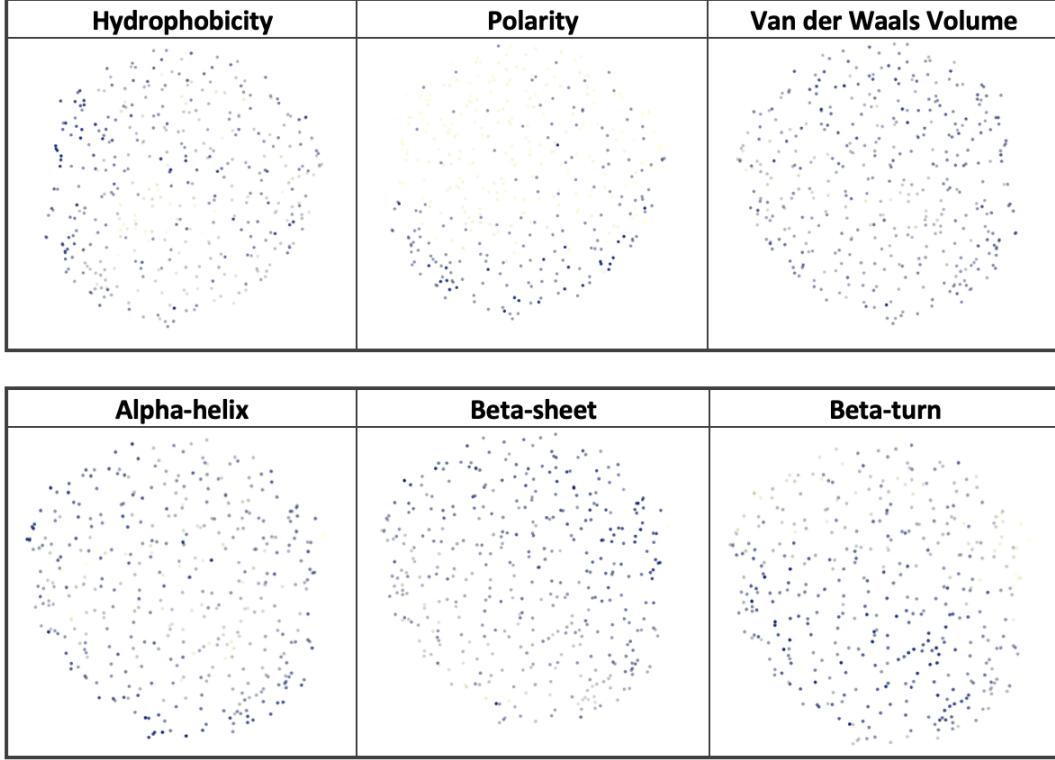


Figure 3: t-SNE 2-gram Embedding

the outer region. As  $\alpha$ -helices have one hydrophobic side and one hydrophilic side, on average half the residues would be hydrophobic and half polar, which is exactly what the embedding displays. Therefore, Figure 3 shows a suitable encoding of biological information for the input representation via an embedding of 2-grams.

#### 4.2.2 Output Representation

Polymers can be mathematically represented by the Cartesian coordinates of their atoms, or by the internal coordinates given by a sequence of bond lengths, angles, and torsions of adjacent bonded atoms. Both parameterisations have their advantages and disadvantages. Spatially proximal residues distant in sequence will be more readily observed in Cartesian space than via their internal coordinates. However, internal coordinates offer a consistent local representation of the structure, whose numerical representation does not deviate from arbitrary global rotations as it would in Cartesian space. Moreover, the internal coordinates for a given protein, folded to its native state, will always be the same.

As the bond lengths and angles are fixed, the proteins can be represented solely via their torsions. Furthermore, 99.9% of residues adopt the trans isomer so the  $\omega$  torsion can be disregarded. This leaves the  $\phi$  and  $\psi$  angles as the only variable components of the internal coordinate system. As a result of this, the output representation consists of discretised tuples of the  $\phi$  and  $\psi$  angles. Discretisation creates a discontinuous approximation of the continuous torsional rotations that facilitates construction of the output vocabulary. For simplicity, and time considerations, the angles are discretised to 1 decimal place of their value in radians. This simplification does introduce a maximum error of 0.05 radians ( $2.86^\circ$ ) in both the  $\phi$  and  $\psi$  rotations of each residue in a sequence. However, it has the added benefit of keeping the input and output vocabularies at a comparable

size. More importantly, the main advantage of using this parameterisation is that the predicted torsions for each residue can only be drawn from a distribution of legal angles. Disallowed regions generally involve steric hindrance between the side chain and main chain atoms. Hence, it will never be the case that a pair of  $\phi$  and  $\psi$  angles are drawn from a blank region of the Ramachandran plot (Figure 4) and hence steric hindrance of adjacent residues will never occur. These regions are sterically disallowed for all amino acids except glycine which is unique in that it lacks a side chain.

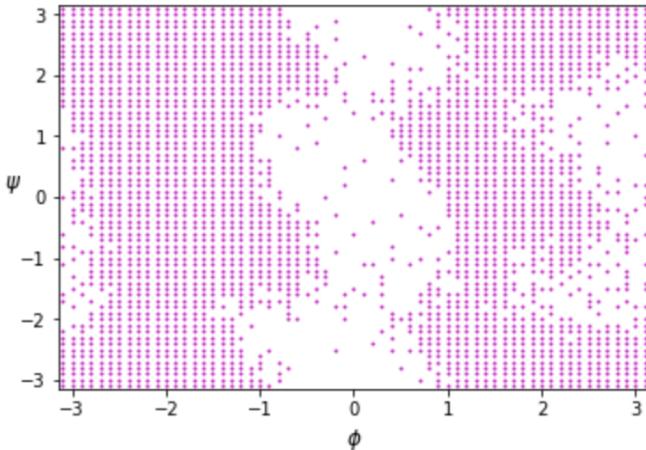


Figure 4: Ramachandran Plot

### 4.3 Data Pipeline

Both the input and output sequences are tokenised using their associated vocabularies, shown in step 5 of Figure 5, whereby each n-gram and torsion pair are mapped to integers from their respective vocabularies. After tokenisation, sequences shorter than the maximum length are padded with zeros in step 6, ensuring that each sequence fed into the network is of the same length. To read data efficiently it is helpful to serialise the data and store it in a set of files. A TFRecord is a simple format for storing a serialised sequence of key-value pairs. For each protein in the dataset, the PDB Code, the tokenized padded amino acid sequence and the tokenized padded torsion sequence are stored.

Achieving peak performance requires an efficient input pipeline that delivers data for the next step before the current step has finished, ensuring optimal use of both the CPU and GPU. Google Colabatory is a free-to-use Jupyter notebook environment that runs entirely in the cloud, allowing free access to a Tesla K80 GPU, with 12GB of RAM, in 12 hour sessions. TensorFlows tf.data API builds efficient input pipelines by framing them as an ETL process:

1. **Extract:** Read data from persistent storage (Google Drive mounted to Google Colab).
2. **Transform:** Use CPU cores to parse and perform pre-processing operations on the data such as shuffling and batching.
3. **Load** Load the transformed data onto the GPU that executes the model.

A naïve synchronous implementation results in the GPU sitting idle as the CPU prepares the data and subsequently the CPU sitting idle as the GPU trains the model. Pipelining overlaps the pre-processing with the model execution of a training step. Data is prefetched through the tf.data API, decoupling the time data is produced from the time it is consumed.

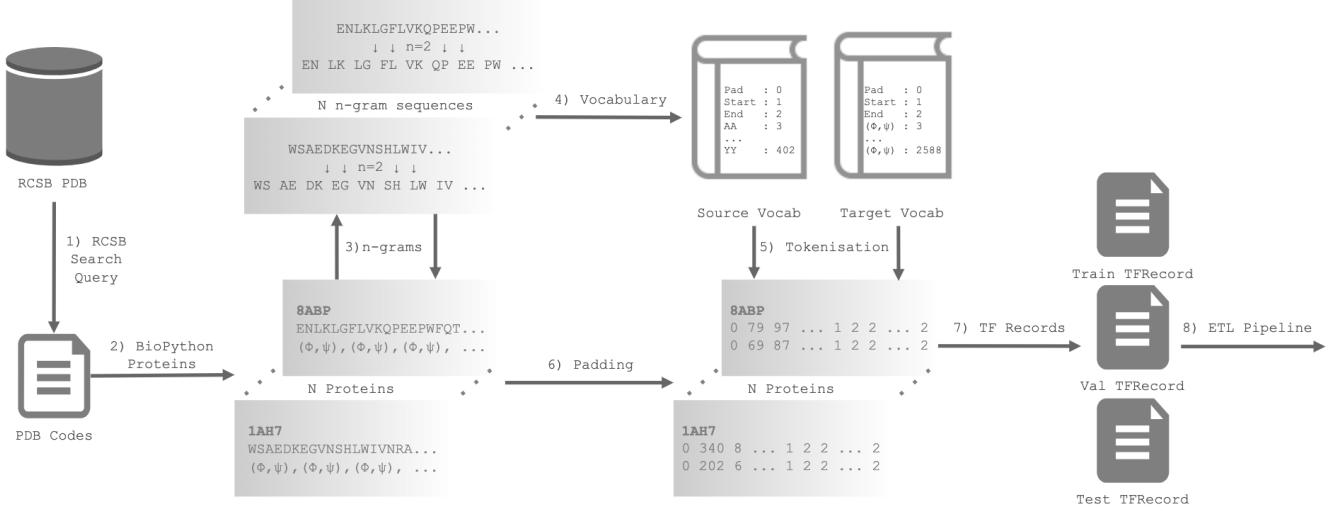


Figure 5: Preprocessing Pipeline

## 4.4 The Network

The initial, rather vanilla, implementation of the neural network model consists of two Recurrent Neural Networks (RNN) that act as an encoder-decoder pair and an attention mechanism [29]. The encoder encodes the variable-length input sequence of tokens into a fixed-length vector representation, and the decoder maps that representation back to a variable-length sequence of output tokens (variable through padding). RNNs are a class of neural network that perform well on sequential data, evident through speech recognition and machine translation tasks [30, 31]. They scan through data from left to right, sharing parameters between each timestep. This ensures that predictions at an arbitrary timestep receive information from all previous states. In order to incorporate information from future states, a Bidirectional RNN (BRNN) was implemented for the encoder [32]. BRNNs have been shown to be significantly more effective than their unidirectional counterparts, especially where context is vitally important [33]. In principle a large enough RNN should be sufficient to generate sequences of arbitrary complexity. In practice however, standard RNNs are unable to store information about past inputs for very long [34]. Hence, both the encoder and decoder RNNs are gated with Long Short-Term Memory (LSTM) units because this RNN architecture is designed to be better at storing and accessing information than standard RNNs [35]. They also help remedy the vanishing gradient problem that plagues conventional Backpropagation Through Time (BPTT).

The intuition behind this design is that the network may be able to learn long-range dependencies between residues distant in sequence, that find themselves close in Cartesian space. Consequently, these dependencies would be analogous to the bonds that form to give proteins their structure. Illustrative examples of this include the long-range dependency that results in two cysteines forming a disulphide bond or the more short-range dependency of hydrogen bonds forming secondary structural motifs. It is hoped that, as the physiochemical information is well-encoded in the input representation, this should be entirely possible.

### 4.4.1 Optimisation

Stochastic gradient-based optimisation is a method of minimising an objective function  $J(\theta)$  parameterised by a models parameters  $\theta \in \mathbf{R}^d$  through updating the parameters in the opposite

direction to the gradient of the objective function  $\nabla_{\theta}J(\theta)$  w.r.t  $\theta$  [36]. The learning rate  $\eta$  determines the size of the steps taken towards some (often local) minima. Stochastic Gradient Descent (SGD) is an iterative method for stochastically approximating gradient descent optimisation that performs a parameter update for each training example  $x$  and label  $y$ :

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta}J(\theta) \quad (4)$$

SGD performs frequent updates with high variance that cause the objective function to fluctuate, enabling it to jump to new, potentially better, local minima. However, this also reduces the likelihood of converging to the global minima as overshooting it is common. For this reason, mini-batch gradient descent is performed to reduce the variance of parameter updates, leading to a more stable convergence. Additionally, TensorFlows highly optimised matrix operations compute gradients w.r.t mini-batches with greater efficiency. The dataset is shuffled before batching to improve the generalisability of the model. Batch sizes of 64 were found to improve efficiency and generalisability most, while operating comfortably within the limits of the hardware. The new update equation is given by Equation 5, where  $b$  is the batch size:

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta}J(\theta; x_{i:i+b}, y_{i:i+b}) \quad (5)$$

Adam is an alternative method for efficient stochastic optimisation [37]. The name Adam is derived from adaptive moment estimation, as the method computes individual adaptive learning rates for different parameters from estimates of first and second moments. Again the aim is to minimise an objective function  $J(\theta)$  w.r.t its parameters  $\theta$ . However, Adam also updates exponential moving averages of both the gradient  $m_t$  and the squared gradient  $v_t$ , where the hyperparameter  $\beta_1, \beta_2 \in [0, 1]$  control the exponential decay rates of these moving averages. The moving averages themselves are estimates of the 1<sup>st</sup> moment and 2<sup>nd</sup> raw moment of the gradient, or the mean and uncentered variance respectively. Thus, the update rule is given by:

$$g_t = \nabla_{\theta}J(\theta) \quad (6)$$

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t \quad (7)$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \quad (8)$$

$$\eta_t = \eta \sqrt{1 - \beta_2^t} / (1 - \beta_1^t) \quad (9)$$

$$\theta_{t+1} = \theta_t - \eta_t \frac{m_t}{\sqrt{v_t + \varepsilon}} \quad (10)$$

As  $m_t$  and  $v_t$  are initialised to vectors of zeros they are biased towards zero, especially initially or when decay rates are small. To counteract this, bias-corrected first and second moment estimates are computed:

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t} \quad (11)$$

$$\theta_{t+1} = \theta_t - \eta_t \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \varepsilon}} \quad (12)$$

The authors propose default values of:  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and  $\varepsilon = 10^{-8}$ . Learning rates that are too small lead to slow convergence while learning rates that are too large hinder convergence, causing the loss function to fluctuate around some minimum and potentially diverge. Figures 6

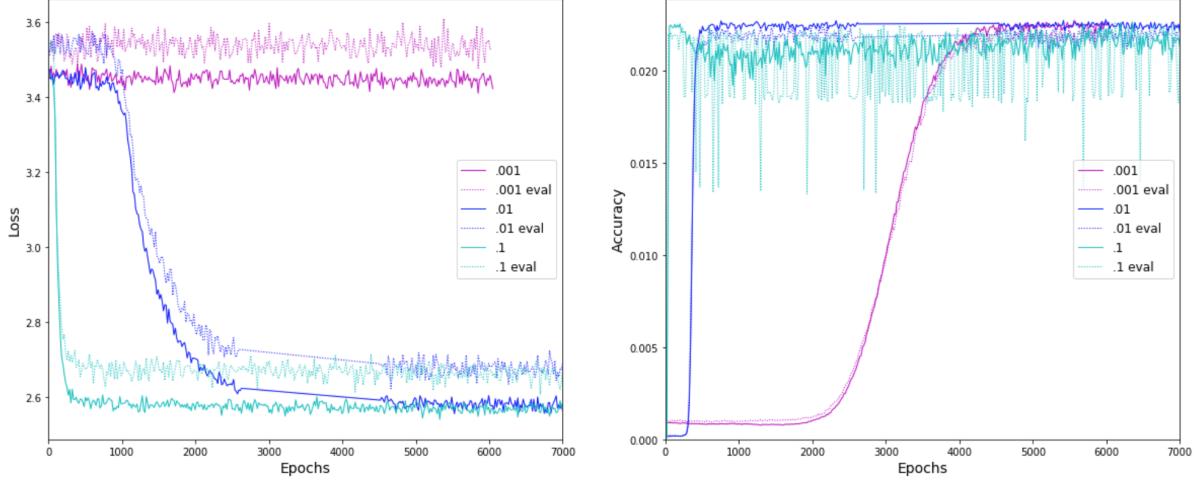


Figure 6: SGD Optimisation

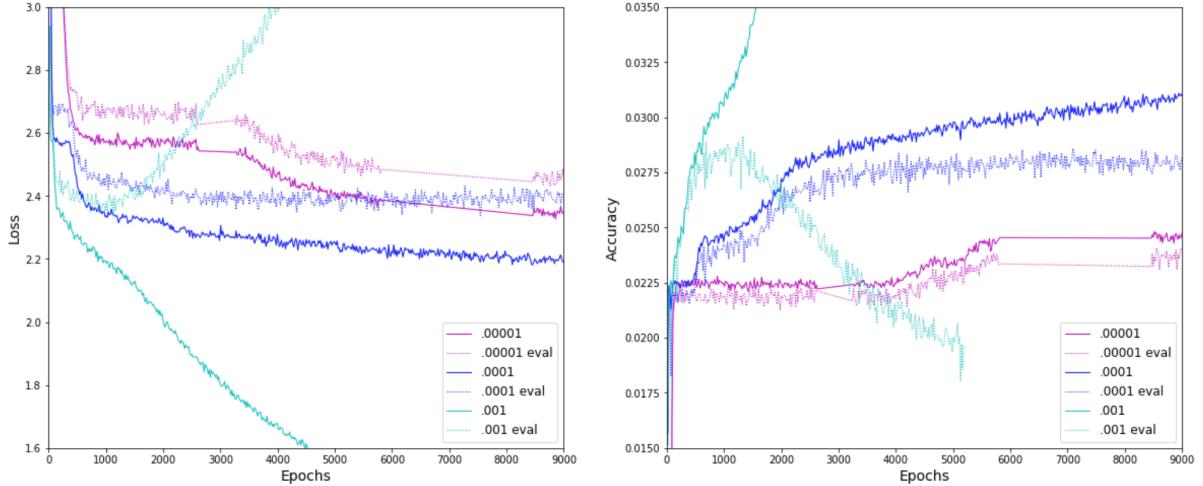


Figure 7: Adam Optimisation

and 7 show experimentation into various learning rates for both SGD and Adam. A comparison between them reveals the more suitable of the two optimisers.

A key challenge of minimising highly non-convex error functions with neural networks is avoiding converging to any of the suboptimal local minima. Moreover, it has been argued by Dauphin *et al.* that non-convex error surfaces in high-dimensional spaces generically suffer from a proliferation of saddle points [38]. It is notoriously hard for SGD to escape saddle points as they are usually surrounded by a plateau of the same error in all directions where the gradient is approximately 0 everywhere. They also state that the ratio between the number of saddle points and the number of local minima increases exponentially with the dimensionality of the function. The results in Figures 6 and 7 are contingent with this insight as SGD struggles to converge while Adam manages to overcome this, converging closer to the optimum without plateauing in either accuracy or loss. As the learning rate of 0.00001 begins to diverge, 0.00005 was chosen as the learning rate for the network from this point on as it falls evenly between the best two values from Figure 7. The Literature Review states that other optimisers such as RMSProp would be analysed, however for time considerations this was omitted as Adam performed best.

#### 4.4.2 Residual Connections

Theoretical and empirical evidence indicates that the depth of neural networks is crucial for their success [39]. With deeper networks, the integration of low-level to high-level features is enriched. Google showed that for NMT systems to achieve good accuracy, both the encoder and decoder have to be deep enough to capture the subtle irregularities in the source and target languages. This decision was made after the success of Sutskever *et al.* from stacking LSTM layers for various NMT problems [40, 41]. In computer vision, models with more than ten convolution layers outperform shallow ones, advancing the state of the art in pattern recognition [42]. The intuition being deeper models possess greater representational power, enabling a hierarchy of concepts to form. However, as depth increases, so too does training time and training difficulty, likely due to exploding/vanishing gradients. Figure 8 displays this well as deeper encoders and deeper decoders converge slower; per both epoch and unit time. Simple stacked LSTM layers in both the encoder and decoder work well up to 4 layers, but 8 layers and beyond saw downgrades in performance.

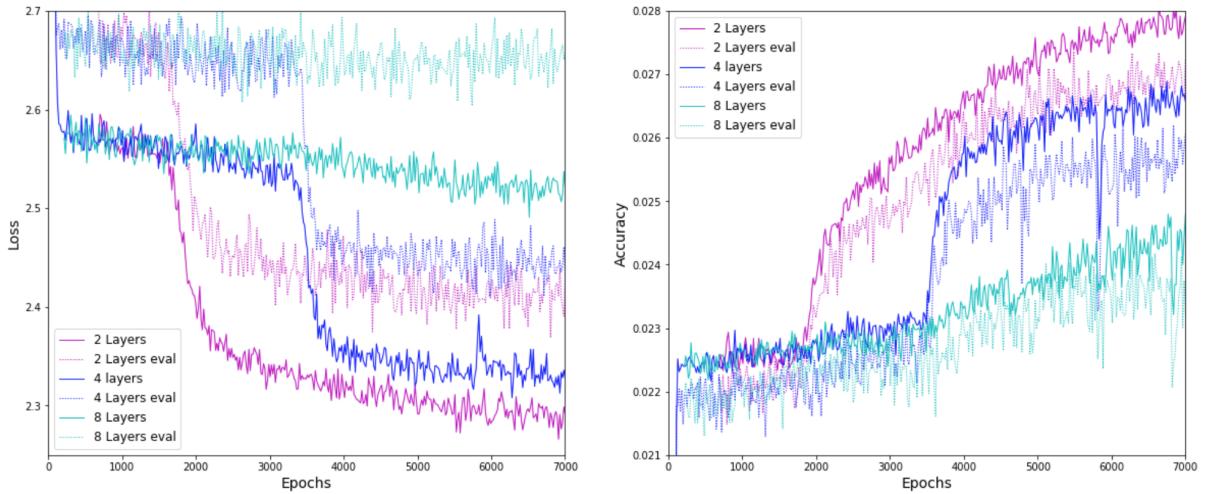


Figure 8: Deeper Networks

The success of simple stacking follows the law of diminishing returns, to the point that by continually increasing the depth, performance can even begin to decline on the training set [43]. As hierarchical composition is, in principle, more powerful, the decline encountered in practice is attributed to the problem of vanishing gradients. The suggestion that adding more layers would decrease representational power is contradictory. Residual networks resolve the problem of training deeper architectures by allowing units to copy their inputs on to the next layer unchanged, via skip connections. For each unit, the output  $y(x)$  is given by Equation 13 where the input  $x$  is summed to a residual  $F(x)$  [44]:

$$y(x) = F(x) + x \quad (13)$$

It is argued that this improves the error flow by decreasing the length of the gradient propagation path [45]. For this reason residual connections were introduced. Input from the lowest LSTM layer is summed element-wise to the layers output and then fed to the subsequent layer.

Training time for 8 residual layers was far too long for the available time of the project to show any reasonable representational benefits from increased depth. Two residual layers did not show much improvement in loss and a decrease in accuracy over the same time-frame. However, 4 residual layers showed improvement in loss but again accuracy began to diverge. As it is expected

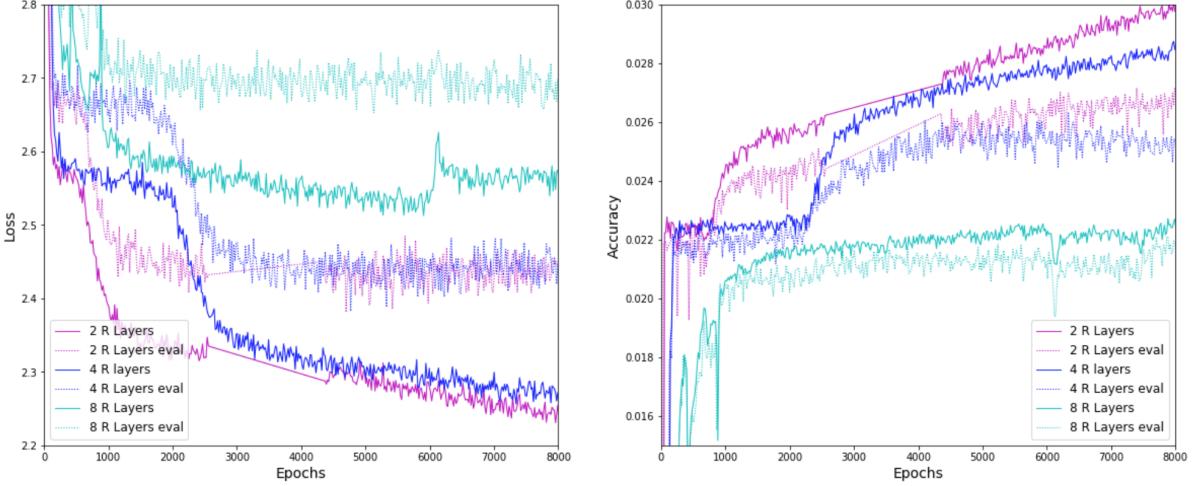


Figure 9: Residual Connections

that deeper representations perform better, a network with 4 residual layers for both the encoder and decoder was chosen, in the hope that divergence in accuracy can be alleviated by the following experiments. Although the Literature Review touched on GNMT, it did not include Residual Connections.

#### 4.4.3 Gradient Clipping

Much like residual connections attempt to resolve a problem attributed to vanishing/exploding gradients, so too does gradient clipping. As introduced by Bengio, the exploding gradient problem refers to the large increase in the norm of the gradient during training [46]. Pascanu *et al.* hypothesised, through a geometrical interpretation, that when gradients explode so too does the curvature along some direction(s), leading to a steep wall in the error surface [47]. Upon a gradient descent step at the cliff edge, an optimiser will be forced to jump across the valley perpendicular to the wall, which could result in a lack of necessary exploration.

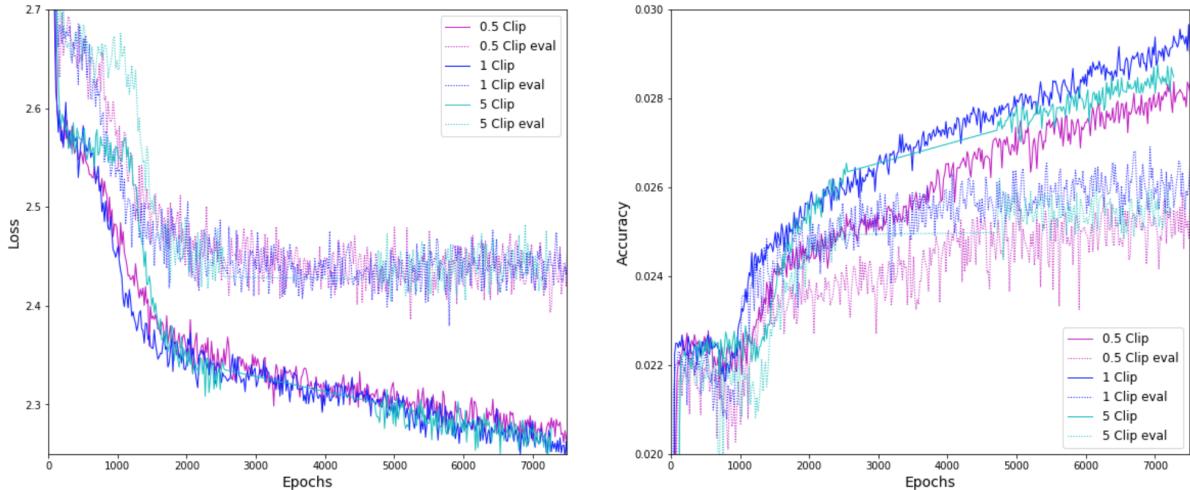


Figure 10: Gradient Clipping

Clipping gradients above a given threshold avoids exploding gradients by trimming the mag-

nitude of the step in the direction that is perpendicular to the edge of the cliff, whilst keeping the direction unchanged; allowing optimisers to more freely explore such valleys. Thus gradient clipping implements a simple form of second-order optimisation in the sense that the second derivative is also proportionally larger in those exploding gradient regions [48]. Gradients are clipped by the threshold values of 0.5, 1 and 5. The results show that the effect of gradient clipping is limited, although subtle improvements in accuracy were realised when gradients were clipped to 1. The previous Literature Review made no mention of gradient clipping so this has been added since the time of it's writing.

#### 4.4.4 Dropout

Overfitting occurs when the neural network is so closely fitted to the training set that it fails to generalise to the validation set. Examples of overfitting exist in all previous experiments where it is clear that the validation loss begins to diverge from the training loss. The dataset was split 60:20:20 during pre-processing, allowing cross-validation of the results.

With limited training data, as in this case where using well-resolved proteins drastically reduces the dataset size, some arbitrary relationships between inputs and outputs learnt by the many non-linear activations within the network may well be the result of sampling noise. The relationships will exist in the training set but not in the general case, and hence display an example of overfitting. Dropout is a regularisation strategy for neural networks whereby units are dropped at random and thus their connections along with them [49]. Dropout acts as an approximate method of combining exponentially many different neural network architectures efficiently [50]. A range of drop values were experimented with to find the best at preventing complex co-adaptions on the training data. The dropout operator is only applied to non-recurrent connections. By not using dropout on the recurrent connections, the LSTM can benefit from dropout regularisation without sacrificing its valuable memorisation ability [51]. The Literature Review did not reference regularisation strategies, so dropout has been added for this report.

Dropout introduced the largest improvement to the accuracy, achieving a peak value of 2.8%.

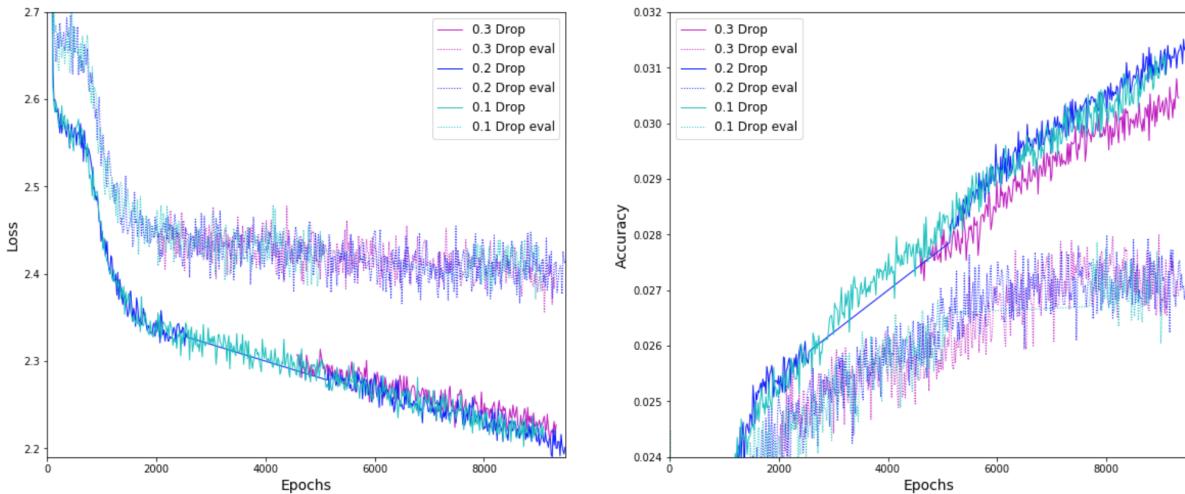


Figure 11: Dropout

## 4.5 Attention

Upon investigation into the limitations of encoder-decoders, Cho *et al.* showed that NMT performs well on relatively short sequences but performance degrades rapidly as the length of the sequence increases [52]. Carrying on from this, Bahdanau *et al.* went on to conjecture that the use of a fixed-length vector is a bottleneck in improving the performance of the encoder decoder architecture [53]. They proposed an attention mechanism that allowed them to search for parts of the source sequence that are relevant to predicting a target symbol, achieving state-of-the-art performance in NMT tasks. Attention has also shown promise in computer vision problems, such as image captioning [54].

Each time the model generates a target symbol it searches for a set of positions in the input sequence where the most relevant information is concentrated. Attentions distinguishing feature is that it no longer encodes the input sequence into a single fixed-length vector. Instead, the input is encoded to a sequence of vectors and chooses a subset adaptively while decoding the output.

The probability is conditioned on a distinct context vector  $C_i$  for each target symbol  $\hat{y}_i$ . The context vector  $C_i$  depends on a sequence of annotations  $(n_1, \dots, n_T)$  where each  $n_i$  contains information about the whole input. The context vector  $C_i$  is computed as a weighted sum of the annotations where  $e_{ij} = a(h_{i-1}, n_j)$  is an alignment model that scores how well the inputs around position  $j$  match outputs at position  $i$ .

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^T \exp(e_{ik})}, \quad C_i = \sum_{j=1}^T \alpha_{ij} n_j \quad (14)$$

Bahdanau successfully applied Equation 14 to jointly translate and align words. Luong drew from this idea to design two novel types of attention based models: a global approach in which all source words are attended and a local one whereby only a subset of source words are considered at a time [55]. The former resembles Bahdanau but the latter is referred to as Luong attention.

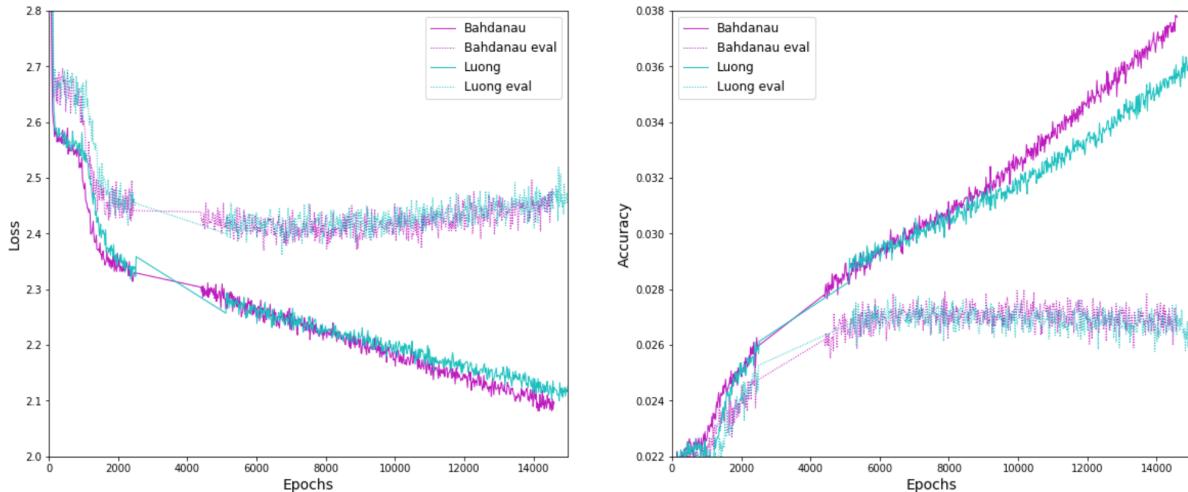


Figure 12: Bahdanau vs. Luong

The global attention of Bahdanau has the drawback that it has to attend all source tokens for each target word. To address that deficiency Luong *et al.* took inspiration from the attempt at image caption generation task by Xu *et al.* that explored the tradeoff between soft and hard attention models. Soft attention refers to global attention, placing a weight over each and every

input. Hard attention, on the other hand, attends to one token at a time. Bahdanau attention was used for all previous experiments. Figure 12 shows that the more global Bahdanau attention performs better at this task than Luong’s attention. However, ungratifying experimentation was carried out in this instance in hope of performance benefits. In the Literature Review it was not specified which attention mechanism would be used but as Bahdanau performs slightly better it was consequently kept as the model’s attention mechanism.

## 5 Implementation

The final implementation of the model consists of a deep LSTM network with 4 encoder layers and 4 decoder layers. The encoder layers are all bidirectional with 128 units in each and all the decoder layers are unidirectional. Residual connections are used for both the encoder and decoder. A Bahdanau attention mechanism connects the lowest decoder layer to the highest encoder layer with an attention size of 128. The Adam optimiser was used for training with a learning rate of 0.00005. Batch sizes of 64 are fed into the network. Dropout is applied to the non-recurrent connections at a drop rate of 0.3. The gradients are clipped at 1.0. Figure 13 shows a diagram of the network.

Early stopping was used after 7000 training steps, providing a model at maximum performance. The minimum loss reached by the final implementation on the validation set is 2.39 with an accuracy of 0.028. The Literature Review intended to use the mean-squared error as the loss metric. However upon further consideration this was seen as a rather naïve metric and so sparse categorical cross entropy was used instead.

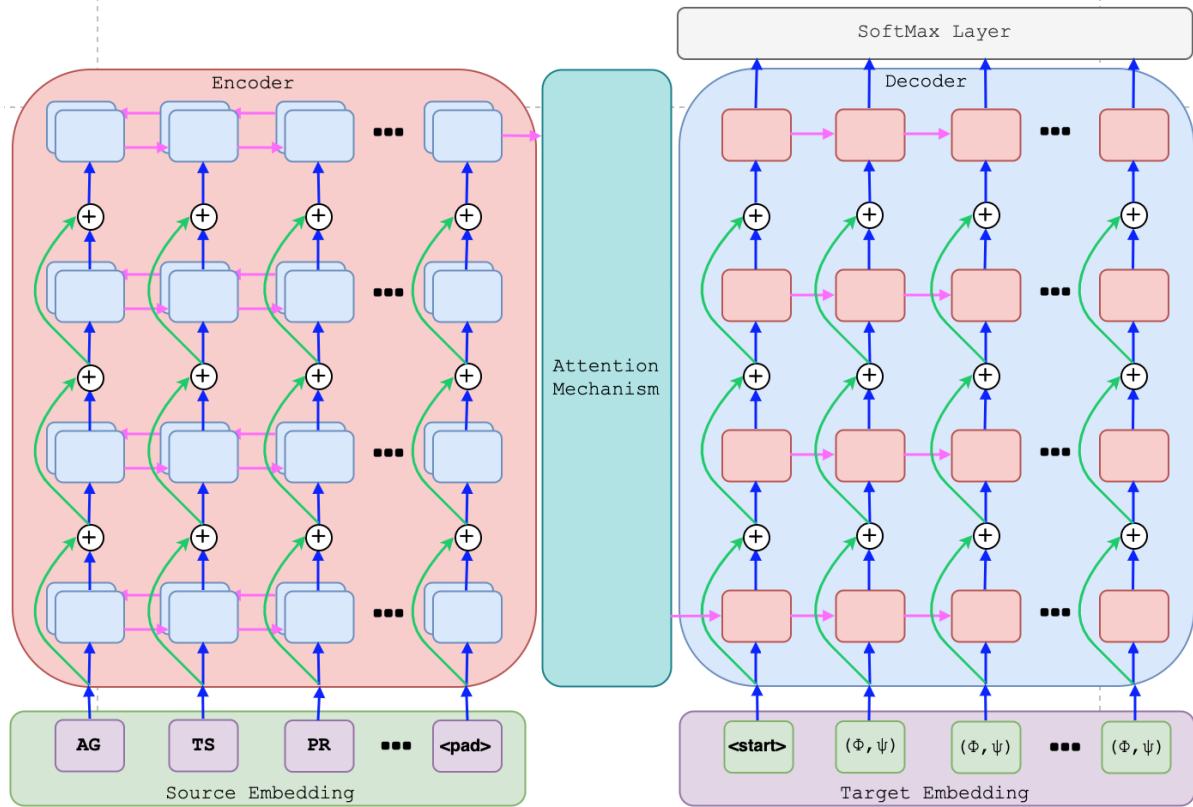


Figure 13: The Model

## 5.1 Inference

The encoder-decoder model learns a conditional distribution over a sequence conditioned on another sequence  $P(\hat{y}_1, \dots, \hat{y}_{T_y} | x_1, \dots, x_{T_x})$ . To find the closest sequence of torsions for a given sequence of n-grams this conditional probability is maximised. This is typically done with a beam search over a greedy search in order to generate an entire sequence of word that maximises the joint probability as opposed to just the most likely word. Therefore, to calculate the probability of a predicted sequence  $\hat{Y}$ , the decoder decomposes the joint probability into ordered conditionals.

$$P(\hat{Y}) = \text{argmax} P(\hat{y}_1, \dots, \hat{y}_{T_y} | x_1, \dots, x_{T_x}) = \prod_{t=1}^T P(\hat{y}_t | y_1, \dots, y_{t-1}, C) \quad (15)$$

The beam width  $B$  specifies the number of torsion to be considered at each step. With larger  $B$  more possibilities can be explored; potentially leading to higher accuracy at the expense of computation time. A beam width of 10 was used.

## 5.2 Torsion Space to Cartesian Space

Efficient conversion from torsion space to Cartesian space is necessary in order to construct PDB files for the predicted proteins. Natural extension Reference Frame (NeRF) is commonly used to switch between parameterisations [56]. Although NeRF can be parallelised to process multiple polymers simultaneously, it is not parallelisable along the length of the polymer. AlQuraishi derived a mathematically equivalent algorithm to overcome this deficiency, pNeRF [57].

A summary of NeRF facilitates understanding of pNeRF. For a given sequence of internal coordinates (encoded via a triplet of bond length, angle and torsion), linear polymers are constructed sequentially by NeRF from one end of the polymer to the other. Coordinates are computed atom by atom in a special reference frame, which, using an affine transformation derived from the three previous coordinates, gives the actual position of the atom. Requiring the previous three atoms to place the next does not allow parallelisation along the length. Thus, pNeRF fragments the polymer into  $M$  segments and converts each to Cartesian space independently before reassembling the fragments with  $M$  affine transformations.

As the model predicts a coarse-grained protein structure, a basic PDB file is created using just the backbone atoms. For each backbone atom in the chain the Cartesian coordinates are written in the PDB file. A quick pass through Rosetta adds the side chains to the file.

## 5.3 Global Distance Test Total Score

Global Distance Test (GDT) Total Score (TS) is a measure of similarity between two protein structures with identical amino acid sequences but varying tertiary structures. This metric is more accurate than RMSD that suffers from sensitivities to outlier regions of a structure that is otherwise reasonably accurate. For example, given an almost perfect match with an unstructured tail comprising a 20Å difference between the two models, this difference will dominate the RMSD, even though the difference may be limited to a small fraction of the structure.

CASP use GDT\_TS measurements as their major assesment criteria, where the higher the GDT\_TS is, the better the prediction. The basic algorithm is to iteratively align two structures under a set of different cutoffs, and then tally the fraction of residues which fall within the cutoffs. The resulting scores end up in the range of 0.0 to 1.0, with 1.0 being a perfect match.

## 6 Results

It is evident from the results that given the limited hardware and training time, the model does not perform well when compared to state of the art methods. AlphaFold averages a GDT\_TS score of 59.65, whereas the implementation of this paper averages a GDT\_TS of around 0.06. It is clear from Figure 14 that improved performance is solely contingent on the frequency of  $\alpha$ -helices within the protein. 2XBG is largely comprised of  $\beta$ -sheets, 2DHO has a balanced mixture of each and 4X2R is mainly  $\alpha$ -helices. Thus, the respective GDT\_TS scores of 0.03, 0.08 and 0.105 are understandable. Every single protein predicted by the model contains an  $\alpha$ -helix down the entire length of the chain. The reason for this is analysed in the following section.

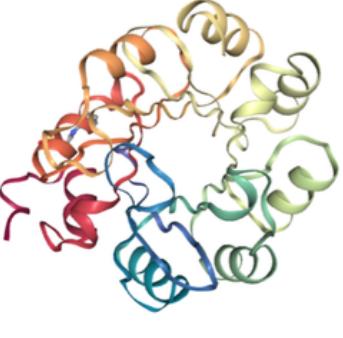
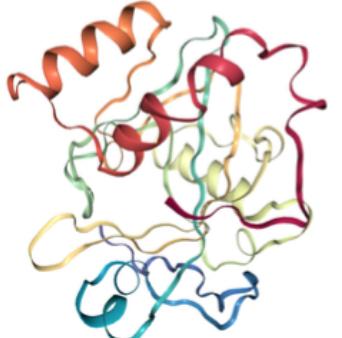
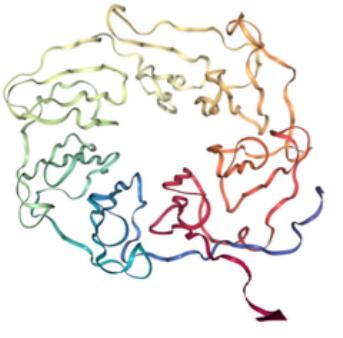
	<b>Best: 4X2R 0.105</b>	<b>Average: 2DHO 0.08</b>	<b>Worst: 2XBG 0.03</b>
<b>Actual Structure</b>			
<b>Predicted Structure</b>			

Figure 14: Results

A scatter plot of the summed difference between the real torsions and predicted torsions against the length of the protein chain for each protein shows a strong correlation between the length of the protein and the accuracy of the prediction. Each point is a protein from the test set. It was only noted upon analysis of the test data that some proteins fell below the minimum sequence length of 200. Despite the search query, which stated only singular chains between 200 and 400 residues should be included, proteins with multiple chains were returned by the query. This meant that some proteins were much smaller than the 200 residue threshold which severely skewed the GDT\_TS results. The smaller the protein, the easier it is to align to a segment of the target protein, especially when it is a simple  $\alpha$ -helix. Figure 14 shows the best, average and worst proteins that fall within the proper range of protein length.

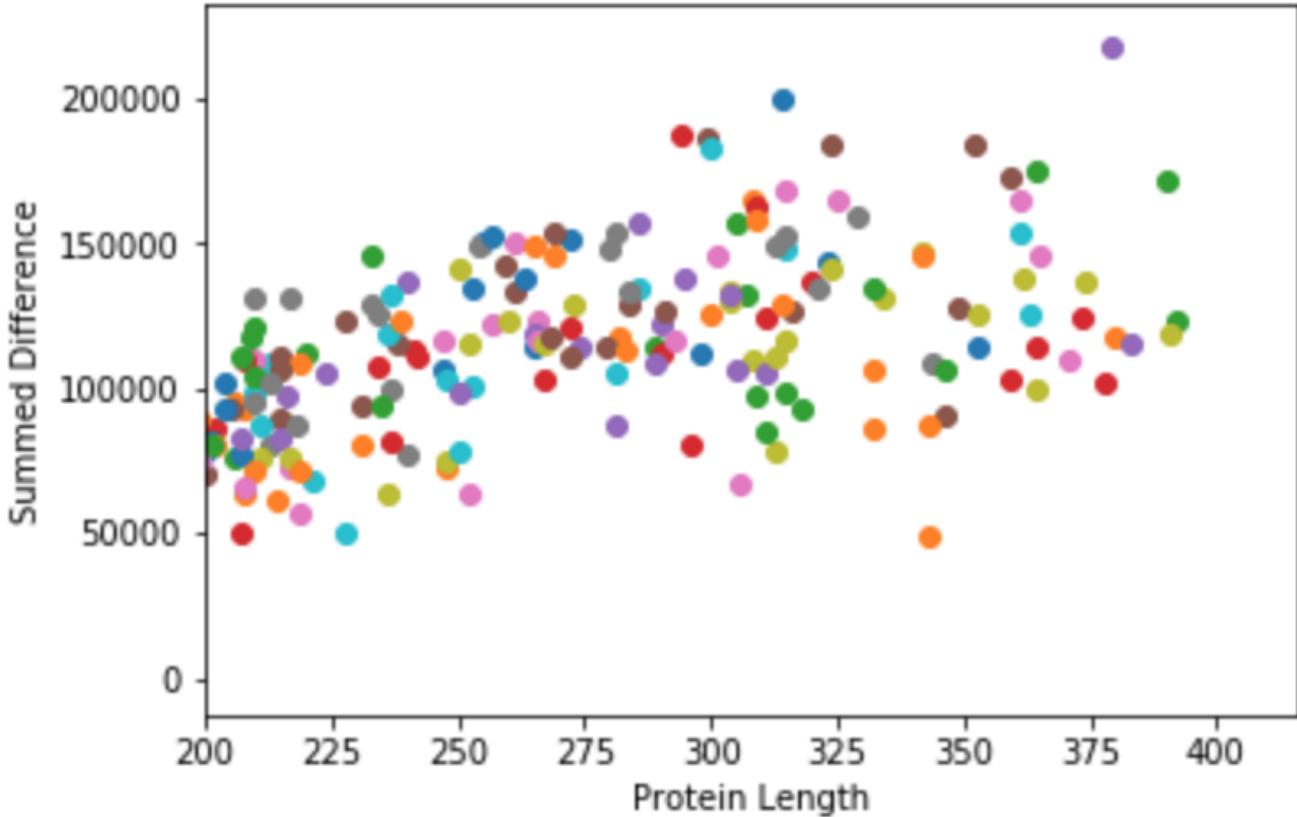


Figure 15: Length Analysis

## 7 Analysis

### 7.1 Angle Distributions

The most frequent  $\phi$  and  $\psi$  tuple is  $(-1.1, -0.7)$  with a total count of 26896 occurrences in the corpus (Figure 16). The output vocabulary consists of 2588 distinct torsional tuples. Thus, for each residue in a sequence there are 2588 possibilities for that residue. The maximum decoder sequence length is 400 because every sequence is padded to this value. Simply choosing the most frequent  $(\phi, \psi)$  pair for each residue results in an accuracy of 2.6% as  $\frac{26896}{400 \cdot 2588} = 0.026$ , which is roughly the maximum accuracy achieved in Section 4. Inspection into the predicted angles from the decoder demonstrates that the majority of predicted tokens come from this region on the Ramachandran plot.

The typical secondary structure from this particular region tends to be  $\alpha$ -helices; as opposed to  $\beta$ -sheets that occupy the top-left region. This is likely the reason that the model predominantly predicts  $\alpha$ -helices over any other structure. The decoder simply learns to place the more frequent torsions along the sequence length as this achieves maximal accuracy without having to learn any difficult dependencies.

The model does manage to decode the initial torsions correctly, such as the start token. Additionally, it manages to learn where the end token is but fails to apply the appropriate padding, instead placing end tokens until the end of the protein. This is promising and means that given more training time, further concepts may be able to be learned by the model.

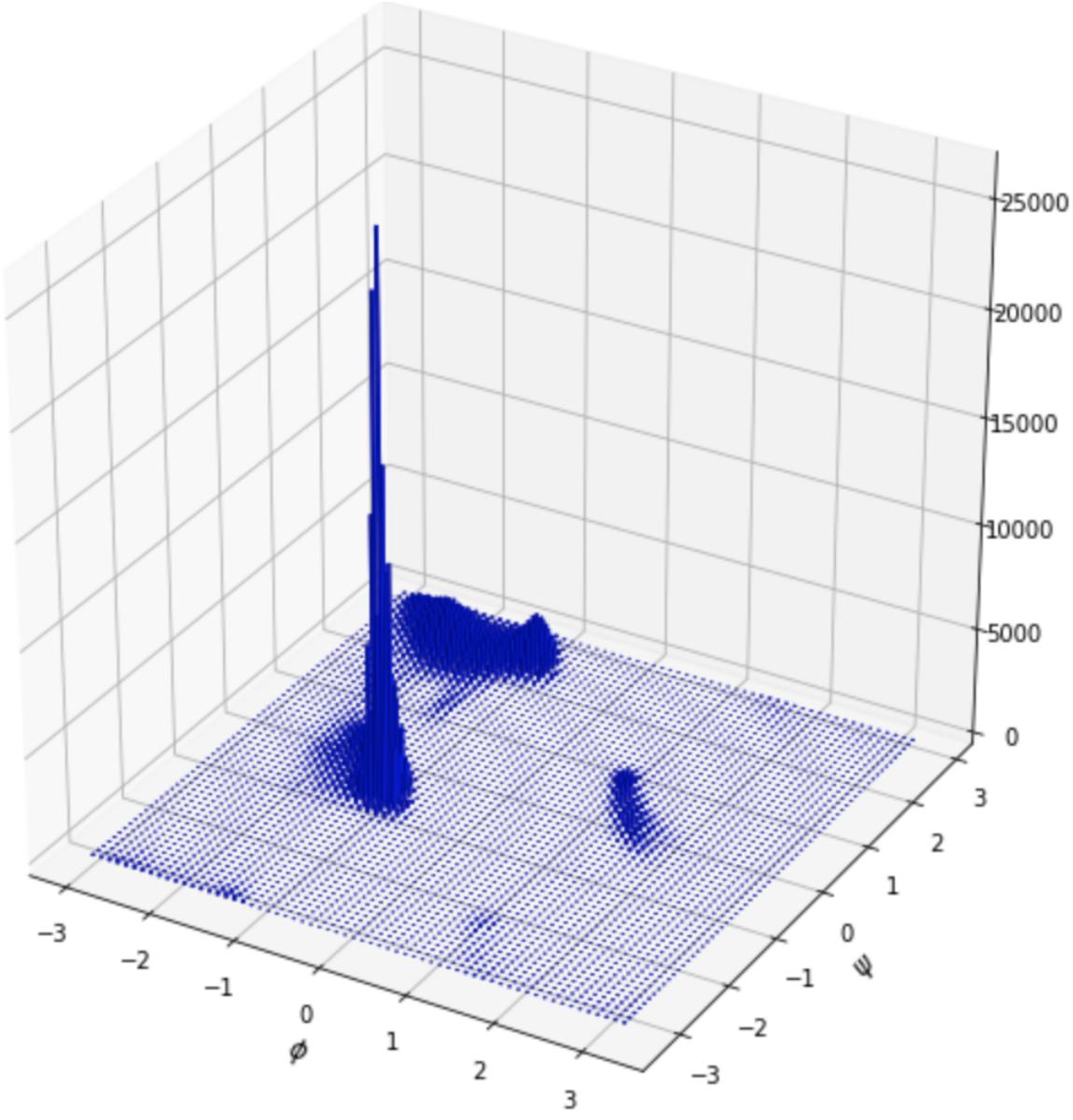


Figure 16: Ramachandran Density Plot

## 7.2 Reconstruction

While assessing the accuracy of pNeRF a detrimental characteristic of the algorithm was encountered. Figure 17 shows Cartesian reconstructions of proteins using pNeRF from their actual internal coordinates, instead of the predicted ones. It is clear that although the  $\phi$  and  $\psi$  angles are correct and secondary structure is conserved, there are various deviations along the chain, resulting in an overall reduction in GDT\_TS even for the true internal coordinates.

After conferring with AlQuraishi as to where this discrepancy came from, it became clear that idealised bond lengths and angles were not being used in the naïve pNeRF algorithm. It was assumed in Section 4.2.2 that the bond lengths and angles were fixed. In reality, these mean values of the bond lengths and angles were obtained by Engh and Huber[58]. More recently however, it has been shown that bond lengths and angles are drawn from distributions unique to each residue. Furthermore, these distributions change subject to secondary structure [59]. The

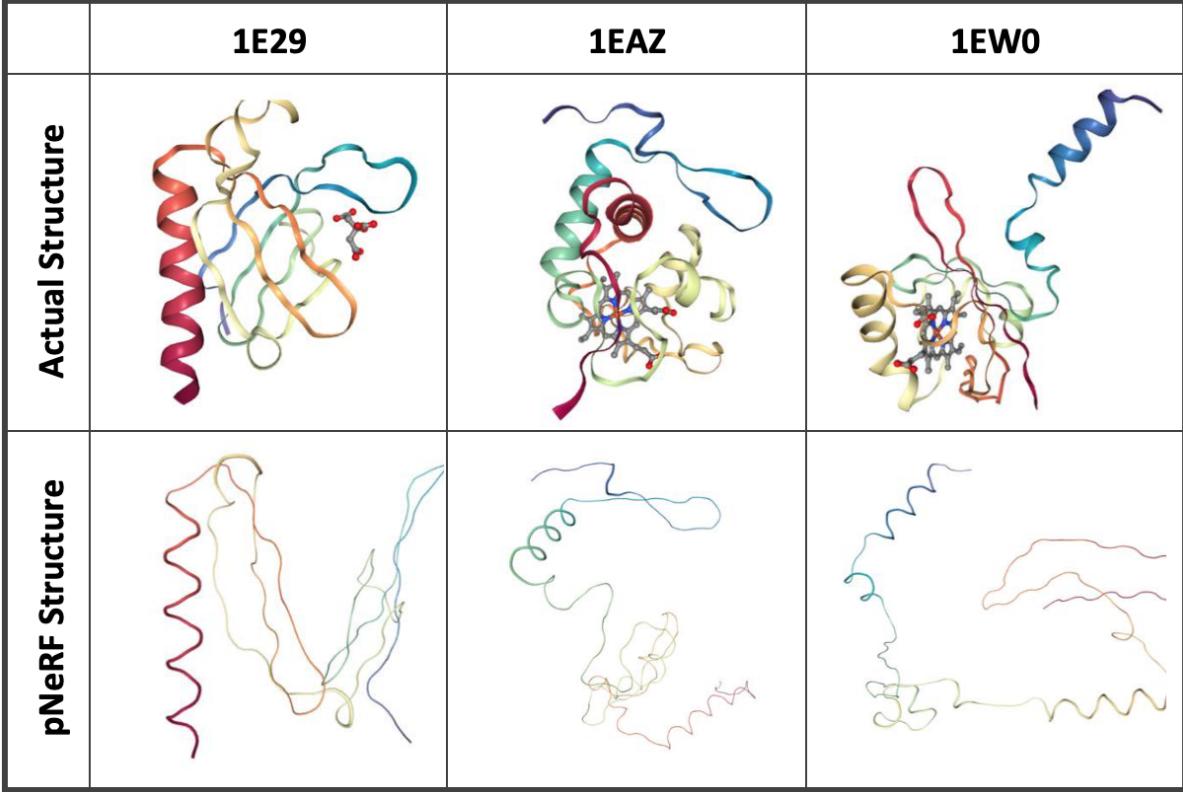


Figure 17: pNeRF Real Reconstruction

pNeRF algorithm simply uses the overall mean bond lengths and angles, regardless of residue and secondary structure. Even if the neural network was as accurate at predicting the internal coordinates as the state of the art methods, this discrepancy in the pNeRF algorithm would limit the performance of the model in terms of the GDT\_TS.

### 7.3 Attentiveness

The x-axis and y-axis of each plot in Figure 18 correspond to the n-grams in the source sequence and the generated internal coordinates, respectively. Each pixel shows the weight  $\alpha_{ij}$  from Equation 14 of the annotation of the  $j$ -th source n-gram for the  $i$ -th target torsion. Unfortunately, the attention mechanism fails to learn any useful information. This is likely due to the length of the proteins. Figure 18 show that the attention is either concentrated in too small a region or spread too thinly over the whole input sequence at each decoding time step. Initial experiments with shorter proteins ranging from 50 residues to 100 residues showed that when decoding a target angle with a particular type of secondary structure, the model could attend to regions in the source sequence of the same secondary structure. However, the size of the dataset for proteins between 50 and 100 residues is limited even further as most proteins are larger than 100 residues in length. A longer training time may be needed in order for the attention mechanism to become fully effective. Alternatively, a larger attention size may offer improvements.

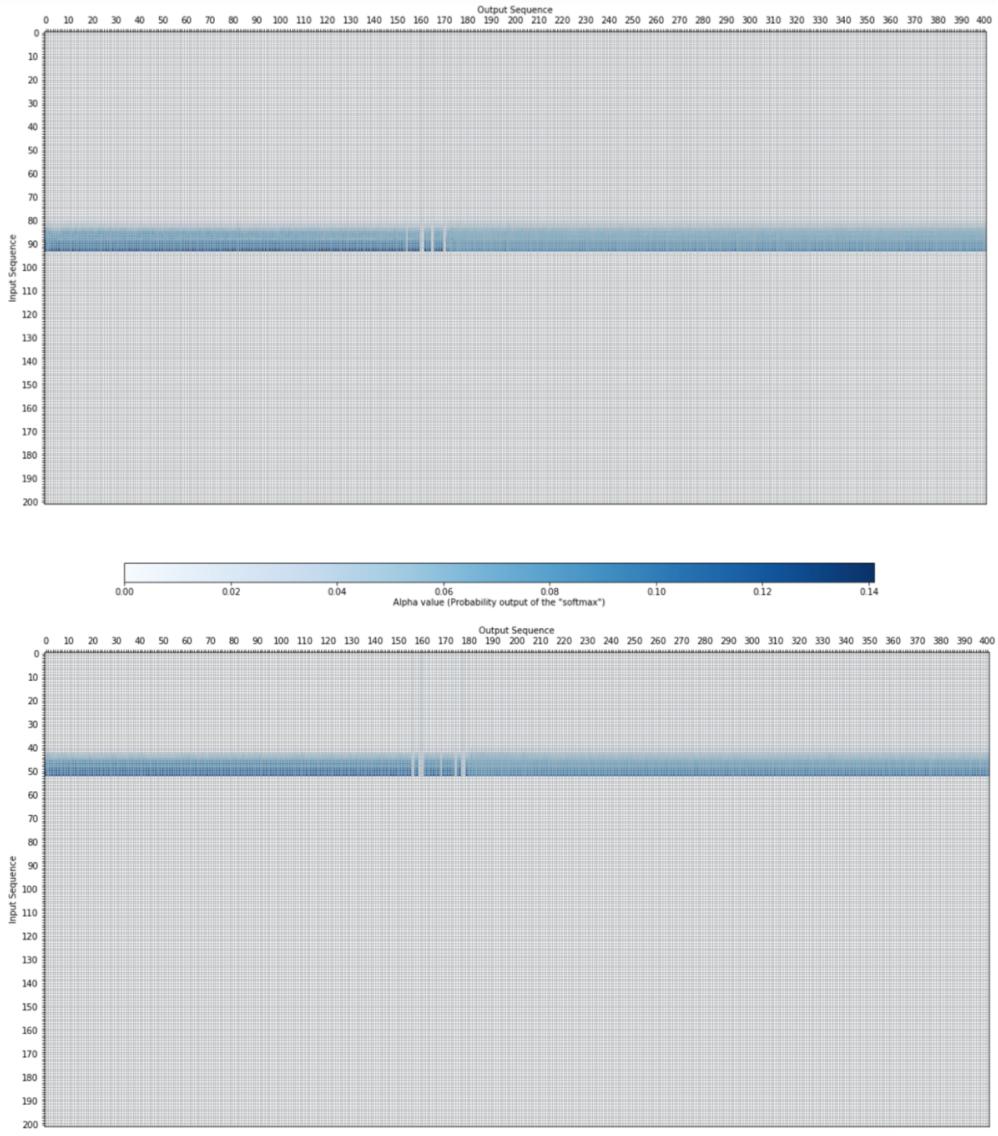


Figure 18: Attention Maps for 41Q6 (top) and 6QPR (bottom)

## 7.4 A Landscape Comparison

Proteins fold in high-dimensional energy landscapes through myriads of confirmations [60]. The energy landscape theory of protein folding is a statistical description of a protein's potential surface [61]. It suggests an overall topography of a rugged funnel, with a bias towards the native structure, as the most realistic shape of a protein's energy landscape. To ensure folding occurs on biologically relevant timescales, the overall gradient of the funnel must be steep when compared to the roughness caused by local minima [62]. However, the high frequency of deep accessible minima still creates a highly non-convex funnel, littered with saddle points.

The native structure of the protein is found at the global optima of its energy landscape, which correspondingly is well-predicted near the global optima of the objective function of the neural network. Whether the similarities cease here is unknown, yet certainly intriguing. Especially in light of research into loss surfaces showing networks that create funnel-like landscapes [63]. Choromanska *et al.* suggested that energy/optimisation landscapes of deep neural networks resemble the energy landscape of zero-temperature Gaussian spin glass [64]. Given that some neural networks

hint at similarity with glassy spin systems from statistical mechanics [65, 66, 67, 68], it is promising then that protein folding has often been modelled using spin-glasses [69, 70, 71].

## 8 Future Work

AlphaFold used the second-order gradient descent optimisation algorithm, L-BFGS. L-BFGS is of the quasi-Newton family of optimisation methods, which approximates the Broyden-Fletcher-Goldfarb-Shanno algorithm using a limited amount of computer memory. It uses an estimation of the inverse Hessian matrix to steer its search through variable space, as opposed to the Jacobian matrix used in standard SGD or Adam optimisation. Newton methods have a damped phase and a quadratically convergent phase, meaning they generally require fewer iterations to converge when compared to stochastic optimisation methods. However, this comes at the trade-off of compute time per iteration. If it holds that the error landscape of the loss function resembles the energy landscape of the protein folding process, it would be interesting to see the effect of using this second-order optimisation method to attempt to better navigate the funnel.

Moreover, upon researching the ways in which gradient descent operates on highly non-convex error landscapes, a deep interest in this area of Machine Learning was formed. An in depth analysis of these high-dimensional error surfaces would be an interesting avenue of future research that may reveal insights that ultimately lead to a better engineered approach to the problem. Great examples of this are: [72, 63, 73]. Personally, visualising the intricacies of neural networks helps build a more intuitive understanding of the complex ways in which they work.

After conferring with AlQuraishi about the drawbacks of pNeRF, it became apparent that an improved version of the algorithm could be easily implemented to include the idealised bond lengths and angles. Such an algorithm, pNeRFect, would simply replace the mean bond lengths and angles of each residue in the chain with draws from their respective distributions. Further, secondary structure predictions could be coupled with the model to better draw from these distributions. A quick Monte Carlo simulation would produce a range of idealised configurations that may ultimately improve the GDT\_TS.

Other regularisation strategies were experimented with, such as Ridge Regression, Lasso Regression and Elastic Net regression. However, due to time constraints no added benefit was realised from these experiments. In future, a more comprehensive analysis of the effect of combinations of  $L_1$  and  $L_2$  regularisation may reveal added improvements. It is unclear at the time of writing how the application of these strategies would affect performance. It would be interesting to test the impact of these regularisation strategies on both recurrent connections and non-recurrent connections.

As the model only predicts  $\alpha$ -helices, regardless of the input sequence, it may be useful to take note of GNMTs approach to rare words. In this case, every torsion pair that is not from the  $\alpha$ -helix region of the Ramachandran plot can be considered rare, compared to the  $\alpha$ -helix torsions. It may be necessary to even up the dataset to include an equal amount of  $\alpha$ -helices,  $\beta$ -sheets and  $\beta$ -turns. Experimentation into the affect of levelling this distribution in the training set is likely to yield performance improvements in the model.

In future attempts, more powerful hardware is a necessity. Utilising Google’s ML Engine would enable access to more GPU power as well the number of GPUs. GNMT places each layer of the encoder and decoder on its own specific GPU, enabling a highly parallelised training process. Additionally models with different parameters could be run in parallel to better search the hyper-parameter space which could introduce performance benefits. Google Colab could not handle more

than 2000 proteins as it was limited to 12GB of RAM and preprocessing uses this memory entirely. Again migration to the cloud opens up increased memory such that it would be possible to include more proteins in future, thus increasing the training, validation and testing dataset sizes.

Finally, increasing the depth of the network as well as the number of hidden units, alongside training over a longer time-frame would likely achieve better results. Improved hardware is crucial for this to occur and any future work would have to capitalise on the added benefits from increasing GPU capacity and number.

## 9 Conclusion

This research devised a novel method of protein structure prediction, taking inspiration from AlphaFold and GNMT, framing the task as an NMT problem. It showed that physiochemical information can be well-encoded in the input representation by embedding n-grams in vector space. It showed that the chosen output representation meant legal structures were more likely to form. The report also includes experimentation into various hyperparameters with the goal of improving performance.

However, the approach of the research failed to materialise to competitive results when benchmarked against the contenders of CASP13. Considering the research was conducted independently, and that CASP predictors have whole teams of researchers behind them, the existence of an actual sequence-to-structure mapping that produces accurate secondary structure is promising, even if it is only one secondary structural motif.

Analysis into the performance of the model was conducted, highlighting distributions of angles on the Ramachandran plot as a limiting factor to performance. When analysing attention it appears that the length of the chains may be another limiting factor in the success of the model. Translating languages will vary rarely take a sentence that is more than 200 words long, if ever. A method of ensuring gradient propagation over longer distances may need to be devised and used in conjunction with the attention mechanism.

The report finishes by discussing the various avenues of future research that would be interesting to explore. From understanding error landscapes and their associations with protein folding funnels to improving the algorithms that facilitated this research, there are many directions this research could be taken in with the aim of improving results.

In conclusion, the model devised does predict structure more efficiently than Molecular Dynamics simulations, yet it fails to reach comparable levels of accuracy to the state-of-the-art methods.

## References

- [1] Y. Yuedong, G. Jianzhao, W. Jihua, H. Rhys, H. Jack, P. Kuldip, and Z. Yaoqi. Sixty-five years of the long march in protein secondary structure prediction: the final stretch? *Briefings in Bioinformatics*, 19(3):482–494, 2018.
- [2] F. Alan. *Structure and Mechanism in Protein Science: Guide to Enzyme Catalysis and Protein Folding*. World Scientific, 1999.
- [3] B. Alberts, A. Johnson, J. Lewis, M. Raff, K. Roberts, and P. Walter. *Molecular Biology of the Cell, 4th edition*. Garland Science, 2002.
- [4] G. N. Ramachandran. Stereochemistry of polypeptide chain configurations. *Journal of Molecular Biology*, 7:95–99, 1963.
- [5] G. N. Ramachandran and V. Sasisekharan. Conformation of polypeptides and proteins. *Advances in Protein Chemistry*, 23:283–438, 1968.
- [6] C. B. Anfinsen. Principles that govern the folding of protein chains. *Science*, 181:223–230, 1973.
- [7] C. B. Anfinsen, R. Redfield, Wa. Choate, J. Page, and W. Carroll. Studies on the gross structure, cross-linkages, and terminal sequences in ribonuclease. *Journal of Biological Chemistry*, 2017:201–210, 1954.
- [8] C. B. Anfinsen, E. Haber, M. Sela, and F. H. White. The kinetics of formation of native ribonuclease during oxidation of the reduced polypeptide chain. *Proceedings of the National Academy of Sciences of the United States of America*, 47(9):13091314, 1961.
- [9] F. H. White. Regeneration of native secondary and tertiary structures by air oxidation of reduced ribonuclease. *The Journal of Biological Chemistry*, 236:1353–1360, 1961.
- [10] I. Andrea and S. Carmelinda. Protein structure determination by X-ray crystallography. *Methods in molecular biology (Clifton, N.J.)*, 452:63–87, 02 2008.
- [11] E. Carpenter, K. Beis, A. Cameron, and S. Iwata. Overcoming the challenges of membrane protein crystallography. *Current opinion in structural biology*, 18(5):581–586, 10 2008.
- [12] D. E. Shaw, M. Deneroff, R. Dror, J. Kuskin, R. Larson, J. Salmon, C. Young, B. Batson, K. Bowers, J. Chao, M. Eastwood, J. Gagliardo, J. Grossman, C. Ho, and D. Ierardi. Anton, a special-purpose machine for molecular dynamics simulation. *Communications of the ACM*, 51:91–97, 2008.
- [13] D. E. Shaw. Millisecond-scale molecular dynamics simulations on Anton. Technical report, D. E. Shaw Research, 2009.
- [14] V. Voelz, G. Bowman, K. Beauchamp, and V. Pande. Molecular simulation of ab initio protein folding for a millisecond folder NTL9(1-39). *Journal of the American Chemical Society*, 132(5):1526–1528, 2010.

- [15] A. Kryshtafovych, Bohdan Monastyrskyy, Krzysztof Fidelis, John Moult, Torsten Schwede, and Anna Tramontano. Evaluation of the template-based modeling in casp12. *Proteins: Structure, Function, and Bioinformatics*, 86(S1):321–334, 2018.
- [16] A. Roy, A. Kucukural, and Y. Zhang. I-TASSER: a unified platform for automated protein structure and function prediction. *Nature Protocols*, 5:725 EP –, 03 2010.
- [17] R. Evans, J. Jumper, J. Kirkpatrick, L. Sifre, T. Green, C. Qin, A. Zidek, A. Nelson, A. Bridgland, H. Penedones, S. Petersen, K. Simonyan, S. Crossan, D. Jones, D. Silver, K. Kavukcuoglu, D. Hassabis, and A. Senior. De novo structure prediction with deeplearning based scoring. In *Thirteenth Critical Assessment of Techniques for Protein Structure Prediction (Abstracts)*, 2018.
- [18] M. Remmert, A. Biegert, and A. Hauser. HHblits: Lightning-fast iterative protein sequence searching by hmm-hmm alignment. *Nature Methods*, 9, 2012.
- [19] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, J. Klingner, A. Shah, M. Johnson, X. Liu, L. Kaiser, S. Gouws, Y. Kato, T. Kudo, H. Kazawa, K. Stevens, G. Kurian, N. Patil, W. Wang, C. Young, J. Smith, J. Riesa, A. Rudnick, O. Vinyals, G. Corrado, M. Hughes, and J. Dean. Google’s Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. *CoRR*, abs/1609.08144, 2016.
- [20] F. Bernstein, T. Koetzle, G. Williams, E. Meyer, M. Brice, J. Rodgers, O. Kennard, T. Shimanouchi, and M. Tasumi. The Protein Data Bank: a computer-based archival file for macromolecular structures. *Journal of molecular biology*, 112(3):535–542, May 1977.
- [21] J.M. Zimmerman, N. Eliezer, and R. Simha. The characterization of amino acid sequences in proteins by statistical methods. *Journal of Theoretical Biology*, 21(2):170 – 201, 1968.
- [22] C. Tanford. Contribution of Hydrophobic Interactions to the Stability of the Globular Conformation of Proteins. *Journal of the American Chemical Society*, 84(22):4240–4247, 1962.
- [23] M. Fauchere, J. and Charton, L. Kier, A. Verloop, and V. Pliska. Amino acid side chain parameters for correlation studies in biology and pharmacology. *International Journal of Peptide and Protein Research*, 32(4):269–278, 1988.
- [24] P. Chou and G. Fasman. *Prediction of the Secondary Structure of Proteins from their Amino Acid Sequence*, pages 45–148. John Wiley & Sons, Ltd, 2006.
- [25] Z.. Harris. Distributional structure. *WORD*, 10(2-3):146–162, 1954.
- [26] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean. Distributed Representations of Words and Phrases and their Compositionality. *CoRR*, abs/1310.4546, 2013.
- [27] E. Asgari and M. Mofrad. ProtVec: A Continuous Distributed Representation of Biological Sequences. *CoRR*, abs/1503.05140, 2015.
- [28] L. Maaten and G. Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9, 2008.

- [29] K. Cho, B. van Merriënboer, Ç. Gülcöhre, F. Bougares, H. Schwenk, and Y. Bengio. Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. *CoRR*, abs/1406.1078, 2014.
- [30] A. Graves, A. Mohamed, and G. Hinton. Speech Recognition with Deep Recurrent Neural Networks. *CoRR*, abs/1303.5778, 2013.
- [31] D. Britz, A. Goldie, M. Luong, and Q. Le. Massive Exploration of Neural Machine Translation Architectures. *CoRR*, abs/1703.03906, 2017.
- [32] M. Schuster and K.K. Paliwal. Bidirectional Recurrent Neural Networks. *Trans. Sig. Proc.*, 45(11):2673–2681, 1997.
- [33] A. Graves and J. Schmidhuber. Framewise phoneme classification with bidirectional LSTM networks. *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, 4:2047–2052, 2005.
- [34] S. Hochreiter, Y. Bengio, P. Frasconi, and J. Schmidhuber. Gradient Flow in Recurrent Nets: the Difficulty of Learning Long-Term Dependencies, 2001.
- [35] S. Hochreiter and J. Schmidhuber. Long Short Term Memory. *Neural Computation*, 9:1735–1780, 1997.
- [36] S. Ruder. An overview of gradient descent optimization algorithms. *CoRR*, abs/1609.04747, 2016.
- [37] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2015.
- [38] Y. Dauphin, R. Pascanu, Ç. Gülcöhre, K. Cho, S. Ganguli, and Y. Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. *CoRR*, abs/1406.2572, 2014.
- [39] R. K. Srivastava, K. Greff, and J. Schmidhuber. Training very deep networks. *CoRR*, abs/1507.06228, 2015.
- [40] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to Sequence Learning with Neural Networks. *CoRR*, abs/1409.3215, 2014.
- [41] T. Luong, I. Sutskever, Q. V. Le, O. Vinyals, and W. Zaremba. Addressing the rare word problem in neural machine translation. *CoRR*, abs/1410.8206, 2014.
- [42] K. He, X. Zhang, S. Ren, and J. Sun. Deep Residual Learning for Image Recognition. *CoRR*, abs/1512.03385, 2015.
- [43] R.K. Srivastava, K. Greff, and J. Schmidhuber. Highway networks. *CoRR*, abs/1505.00387, 2015.
- [44] K. Greff, R.K. Srivastava, and J. Schmidhuber. Highway and residual networks learn unrolled iterative estimation. *CoRR*, abs/1612.07771, 2016.

- [45] J. Zhou, Y. Cao, X. Wang, P. Li, and W. Xu. Deep recurrent models with fast-forward connections for neural machine translation. *CoRR*, abs/1606.04199, 2016.
- [46] Y. Bengio, P. Simard, and P. Frasconi. Learning Long-term Dependencies with Gradient Descent is Difficult. *Trans. Neur. Netw.*, 5(2):157–166, March 1994.
- [47] R. Pascanu, T. Mikolov, and Y. Bengio. Understanding the exploding gradient problem. *CoRR*, abs/1211.5063, 2012.
- [48] Y. Bengio, N. Boulanger-Lewandowski, and R. Pascanu. Advances in optimizing recurrent networks. *CoRR*, abs/1212.0901, 2012.
- [49] G. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *CoRR*, abs/1207.0580, 2012.
- [50] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014.
- [51] W. Zaremba, I. Sutskever, and O. Vinyals. Recurrent neural network regularization. *CoRR*, abs/1409.2329, 2014.
- [52] K. Cho, B. van Merriënboer, D. Bahdanau, and Y. Bengio. On the Properties of Neural Machine Translation: Encoder-Decoder Approaches. *CoRR*, abs/1409.1259, 2014.
- [53] D. Bahdanau, K. Cho, and Y. Bengio. Neural Machine Translation by Jointly Learning to Align and Translate. *CoRR*, abs/1409.0473, 2014.
- [54] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhutdinov, R. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. *CoRR*, abs/1502.03044, 2015.
- [55] M. Luong, H. Pham, and C. Manning. Effective Approaches to Attention-based Neural Machine Translation. *CoRR*, abs/1508.04025, 2015.
- [56] J. Parsons, B. Holmes, M. Rojas, J. Tsai, and C. Strauss. Practical conversion from torsion space to cartesian space for in silico protein synthesis. *Journal of Computational Chemistry*, 26(10):1063–1068, 2005.
- [57] M. AlQuraishi. pNeRF: Parallelized Conversion from Internal to Cartesian Coordinates. *bioRxiv*, 2018.
- [58] R. A. Engh and R. Huber. Accurate bond and angle parameters for X-ray protein structure refinement. *Acta Crystallographica Section A*, 47(4):392–400, Jul 1991.
- [59] Wouter G. Touw and Gert Vriend. On the complexity of engh and huber refinement restraints: the angle as example. In *Acta crystallographica. Section D, Biological crystallography*, 2010.
- [60] N. Hori, G. Chikenji, S. Berry, and S. Takada. Folding energy landscape and network dynamics of small globular proteins. *Proceedings of the National Academy of Sciences*, 106(1):73–78, 2009.

- [61] J. Onuchic, Z. Luthey-Schulten, and P. Wolynes. THEORY OF PROTEIN FOLDING: The Energy Landscape Perspective. *Annual Review of Physical Chemistry*, 48(1):545–600, 1997. PMID: 9348663.
- [62] J. Wang, R. Oliveira, X. Chu, P. Whitford, J. Chahine, W. Han, E. Wang, J. Onuchic, and V. Leite. Topography of funneled landscapes determines the thermodynamics and kinetics of protein folding. *Proceedings of the National Academy of Sciences*, 109(39):15763–15768, 2012.
- [63] H. Li, Z. Xu, G. Taylor, and T. Goldstein. Visualizing the Loss Landscape of Neural Nets, 2018.
- [64] A. Choromanska, M. Henaff, M. Mathieu, G. Ben Arous, and Y. LeCun. The loss surface of multilayer networks. *CoRR*, abs/1412.0233, 2014.
- [65] C. Martin and M. Mahoney. Implicit Self-Regularization in Deep Neural Networks: Evidence from Random Matrix Theory and Implications for Learning. *CoRR*, abs/1810.01075, 2018.
- [66] L. Sagun, V. Ugur Güney, and Y. LeCun. Explorations on high dimensional landscapes. *CoRR*, abs/1412.6615, 2014.
- [67] A. Barra, G. Genovese, F. Guerra, and D. Tantari. How glassy are neural networks? *Journal of Statistical Mechanics: Theory and Experiment*, 2012(07):P07009, jul 2012.
- [68] M. Baity-Jesi, L. Sagun, M. Geiger, S. Spigler, G. Ben Arous, C. Cammarota, Y. LeCun, M. Wyart, and G. Biroli. Comparing dynamics: Deep neural networks versus glassy systems. In *ICML*, 2018.
- [69] J. Bryngelson and P. Wolynes. Spin glasses and the statistical mechanics of protein folding. *Proceedings of the National Academy of Sciences*, 84(21):7524–7528, 1987.
- [70] Piotr Garstecki, Trinh Xuan Hoang, and Marek Cieplak. Energy landscapes, supergraphs, and “folding funnels” in spin systems. *Phys. Rev. E*, 60:3219–3226, Sep 1999.
- [71] V.S. Pande, A.Y. Grosberg, and T. Tanaka. Statistical mechanics of simple models of protein folding and design. *Biophysical Journal*, 73(6):3192 – 3210, 1997.
- [72] S. Carter, Z. Armstrong, L. Schubert, I. Johnson, and C. Olah. Activation atlas. *Distill*, 2019. <https://distill.pub/2019/activation-atlas>.
- [73] D. Im, M. Tao, and K. Branson. An empirical analysis of deep network loss surfaces. *CoRR*, abs/1612.04010, 2016.

A PDB Codes

## B Embeddings

### B.1 1-Grams

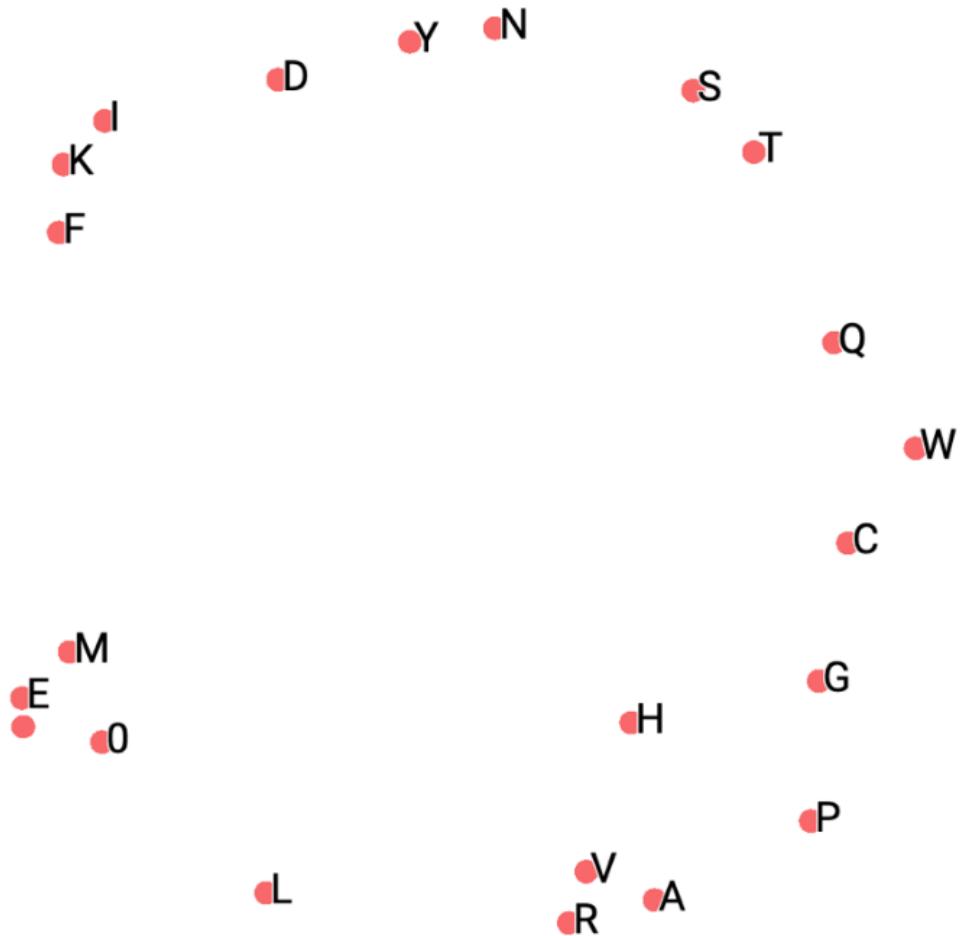


Figure 1: 1-gram

## B.2 2-Grams

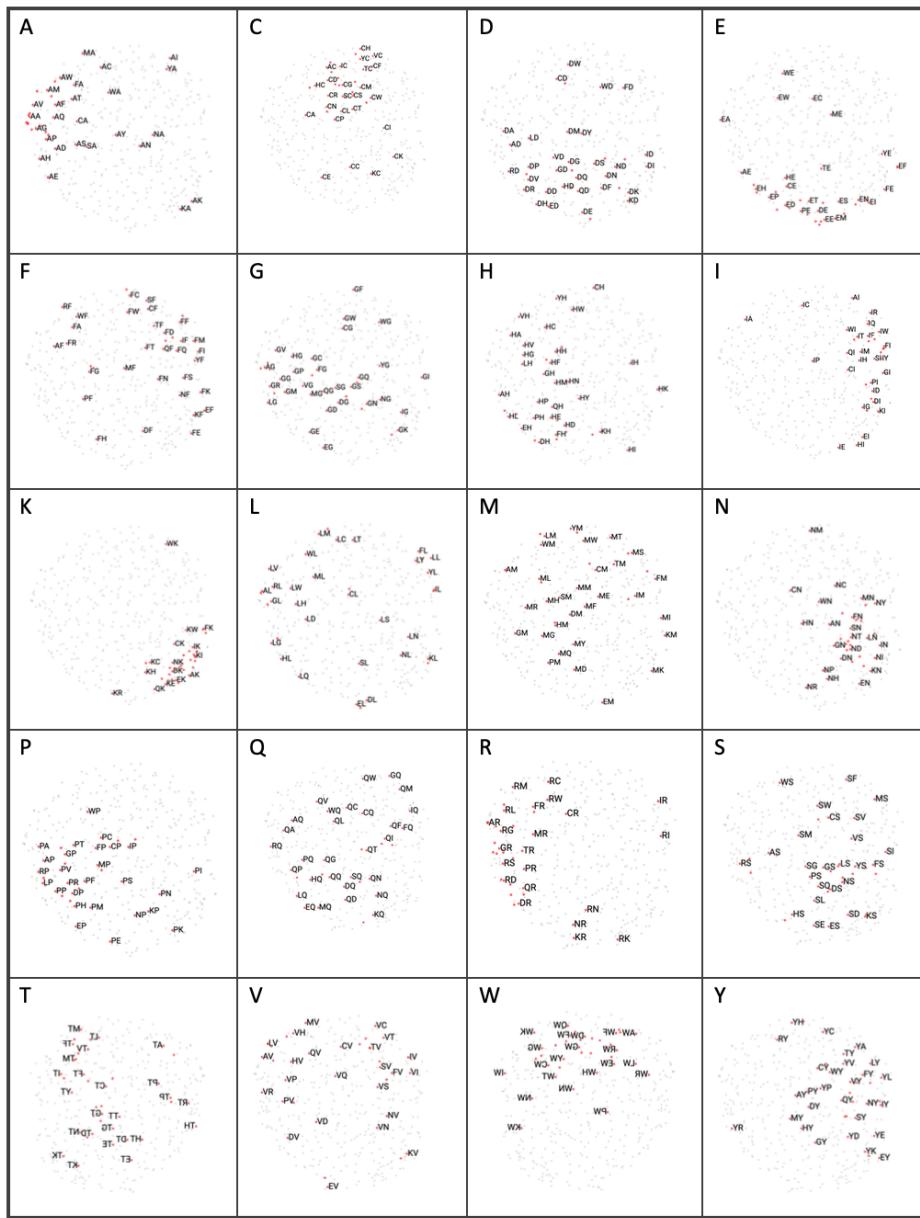


Figure 2: 2-gram

### B.3 3-Grams

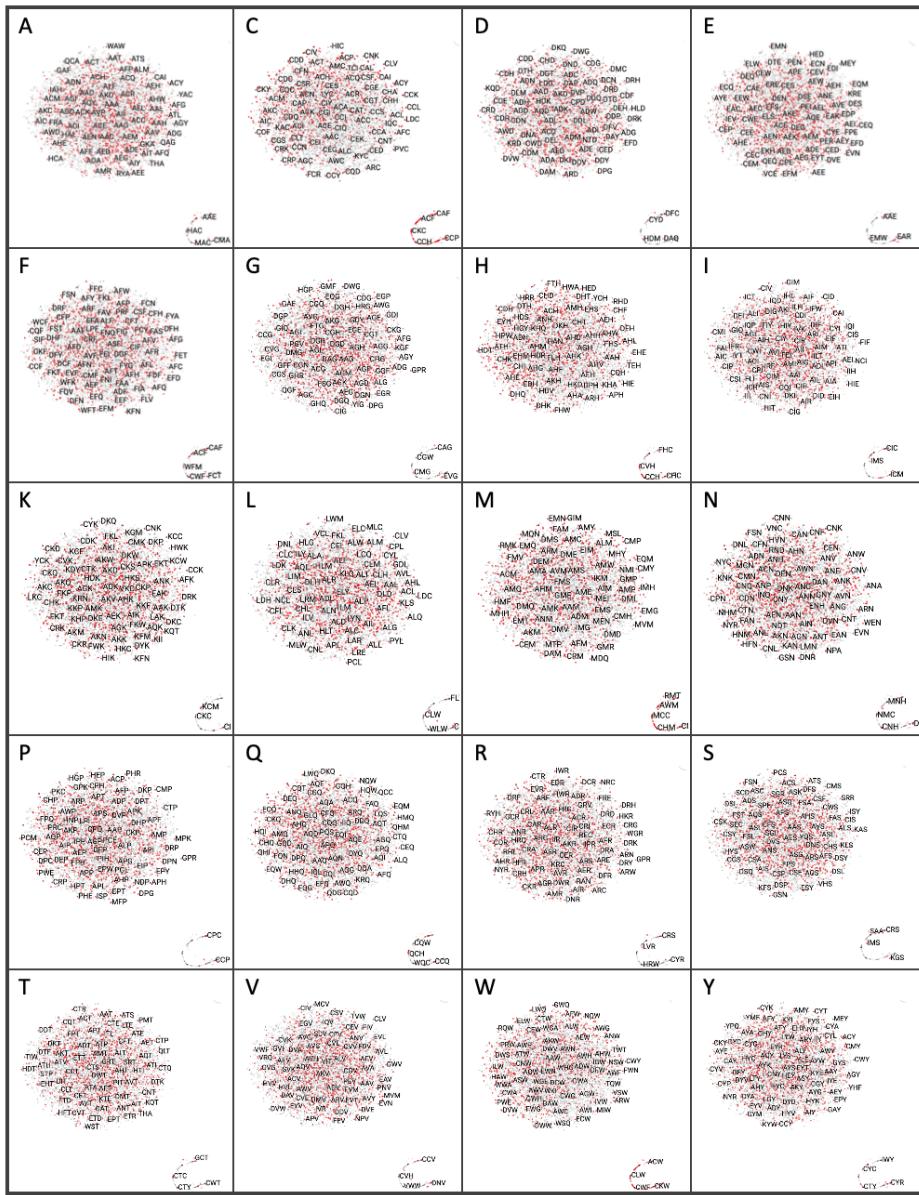


Figure 3: 3-gram