

**Heidelberg University
Institute of Computer Science**

**Project report for the lecture
Advanced Machine Learning**

**Prof. Dr. Köthe
Summer semester 2021**

**Prediction of the next SARS-CoV-2 mutations
by a Transformer based GAN Framework**

<https://github.com/nilskre/AML-covid-project>

Team Member: Felix Hausberger, 3661293,
Applied Computer Science
eb260@stud.uni-heidelberg.de

Team Member: Nils Krehl, 3664130,
Applied Computer Science
pu268@stud.uni-heidelberg.de

Plagiarism statement

We certify that this report is our own work, based on our personal study and/or research and that we have acknowledged all material and sources used in its preparation, whether they be books, articles, reports, lecture notes, and any other kind of document, electronic or personal communication.

We also certify that this report has not previously been submitted for assessment in any other unit, except where specific permission has been granted from all unit coordinators involved, or at any other time in this unit, and that we have not copied in part or whole or otherwise plagiarized the work of other students and/or persons.

Member contributions

Nils Krehl

Came up with the project idea and was the domain expert for the project's biomedical domain. Initial research for the topic and feasibility analysis. Extensive research in the areas of the project's overall composition, related work on mutation prediction, molecular biological basics, data sources, dataset creation and dataset insights. Responsible for the data selection and download, the dataset creation and the generation of dataset insights.

Felix Hausberger

Extensive research in the area of machine learning models, especially the connection between Natural Language Processing (NLP) methods and our project, previous work on mutation prediction and different model architectures such as Sequence-to-Sequence (Seq2Seq), Transformer and Generative Adversarial Networks (GANs). Responsible for the data preprocessing, the model architecture, implementation and configuration, the GAN training and evaluation.

Contents

0 Project Setup	1
1 Introduction	1
2 Fundamentals and Related Work	2
2.1 From Probabilistic Language Models to modeling Evolution Theory	2
2.2 GISAID EpiFlu Data Platform	4
2.3 Domain-Specific Methodologies to create Evolutionary Datasets for Mutation Prediction	4
2.4 Previous Work on Mutation Prediction	6
2.5 Sequence-to-Sequence Models based on Long Short-Term Memory	6
2.6 Applying Generative Adversarial Networks	8
2.7 Transformer and Attention Mechanism	10
3 Approach	14
3.1 Dataset Creation	14
3.1.1 Raw Data Selection from Global Initiative on Sharing All Influenza Data (GISAID)	14
3.1.2 Generation of a Phylogenetic Tree	14
3.1.3 Phylogenetic Tree to Dataset	15
3.2 Data Preprocessing	16
3.2.1 Dimensionality Reduction by selecting Subpart of the Genome	16
3.2.2 Transformation of Genome Sequence to Numeric Model Input	16
3.3 Model Architecture	17
3.4 Training Process	19
4 Experimental Results	21
4.1 Dataset Insights	21
4.1.1 Complete Dataset	21
4.1.2 Preprocessed Dataset	22
4.2 Evaluation	23
4.2.1 Evaluation Criteria	23
4.2.2 Training Phase and Evaluation	24
5 Conclusion	26

A Appendix	I
A.1 Appendix A: Basics from Molecular Biology	I
A.1.1 Genome Sequences	I
A.1.2 Protein Biosynthesis	I
A.1.3 Structure of Severe Acute Respiratory Syndrome Corona Virus 2 (SARS-CoV-2)	III
A.1.4 Relationships between Biological Objects	IV

List of Abbreviations

BLEU	Bilingual Evaluation Understudy
DNA	Deoxyribonucleic Acid
GAN	Generative Adversarial Network
GISAID	Global Initiative on Sharing All Influenza Data
GPU	Graphics Processing Unit
LSTM	Long Short-Term Memory
NLP	Natural Language Processing
NMT	Neural Machine Translation
ORF	Open Reading Frames
ReLU	Rectified Linear Unit
RKI	Robert Koch Institut
RNA	Ribonucleic Acid
RNN	Recurrent Neural Network
SARS-CoV-2	Severe Acute Respiratory Syndrome Corona Virus 2
Seq2Seq	Sequence-to-Sequence
UTR	Untranslated Region

0 Project Setup

For a detailed description of how to set up the project, please have a look at <https://github.com/nilskre/AML-covid-project/blob/main/README.md#setup>.

1 Introduction

During the genome replication, random mutations can appear. As a consequence the encoded protein sequence could be changed, which can lead to different behavior. If this change increases fitness, it is probably passed on to the next generation. For an explanation of the needed basics from molecular biology refer to chapter A.1. [6]

A currently well-known example of a mutating virus is SARS-CoV-2. Due to the developed vaccines, the hope for an end of the pandemic arises. Nevertheless, this is only true, if the vaccines, which are developed against the wild type of SARS-CoV-2, also remain effective against new mutations. To enable fast responses to new arising mutations it would be helpful to know the possible next mutations in advance. This can influence the treatment and prevention of diseases, by enabling the development of countermeasures and preventive measures in advance. [6]

Machine Learning, especially Deep Learning enabled improvements in lots of different domains. This work applies Deep Learning to the area of virus genome mutation prediction. Due to the fact, that genome sequences can be represented as text data using the FASTA format, methods from the NLP area can be applied. The success of Deep Learning for NLP tasks has already been shown in various areas such as text generation, text summarization or translation. [6]

The central research question in this paper is whether a machine learning model can be trained to predict the next possible SARS-CoV-2 mutations. In this paper, we propose three novelties:

- **Model architecture:** A new GAN based architecture influenced by [6]. Central novelty is the usage of transformers instead of an Long Short-Term Memory (LSTM) in the Seq2Seq model.
- **Dataset:** Generation of a dataset for SARS-CoV-2, consisting of 9199 parent-child data instances.
- **Application domain:** The adaption of the model to SARS-CoV-2.

2 Fundamentals and Related Work

2.1 From Probabilistic Language Models to modeling Evolution Theory

A probabilistic language model tries to approximate the probability distribution

$$P(w_1, \dots, w_n) = \prod_{t=1}^n P(w_t | w_1, \dots, w_{t-1}) \quad (1)$$

with w_t being a word at position (time step) t in a sentence of length n . To build language models Recurrent Neural Networks (RNNs) were used to model such probability distributions.

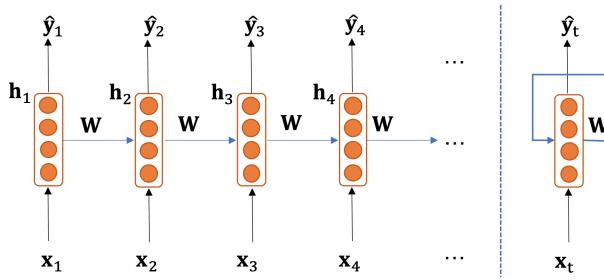


Figure 1: Architecture of a conventional RNN [11, p. 46]

At each time step t it outputs a probability distribution $P(w_t | w_1, \dots, w_{t-1})$ given the words read so far in the current instance (see Figure 1). Words are read as a vectorized numerical representation, often given by pretrained so-called word embeddings x_t which are lower-dimensional and more semantically enriched compared to simple one-hot encodings. One then calculates the hidden state h_t by

$$h_t = f(W^{(h)}h_{t-1} + W^{(x)}x_t + b_1) \quad (2)$$

and the corresponding output probability distribution by

$$\hat{y}_t = \text{softmax}(U^{(h)}h_t + b_2). \quad (3)$$

The applied weight matrix is always the same for each time step t giving the RNN its name. One can therefore simplify the unrolled RNN architecture on the left side of Figure 1 to the one on the right, where the hidden state is continuously passed as an input to the next time step. To achieve a better convergence behavior during training, one can also provide the expected hidden state of time step $t - 1$ instead of using the predicted hidden state,

which is called teacher forcing. RNNs can process input of arbitrary length and are by their recurrent character capable to use information from previous time steps. Unfortunately, they are vulnerable to vanishing and exploding gradient problems. LSTM is a special RNN architecture that solves such vulnerabilities by owning a separate long-term cell state besides a short-term hidden state (introduced in subsection 2.5). It can preserve information over many time steps. [11]

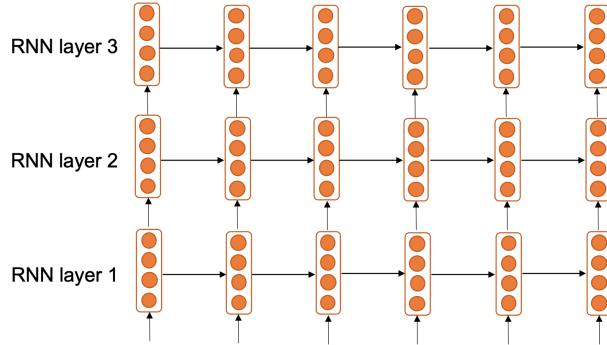


Figure 2: Architecture of a multi-layer RNN [11, p. 87]

One can also use two RNNs, one traversing a sentence from left to right and another one vice versa, with two different weight matrices to model the probability distribution bidirectionally. One therefore simply concatenates the hidden states of each RNN before applying the weight matrix U and the *softmax()* function. Also multi-layer RNN can be utilized to generate higher-order features (hidden states) for the prediction task (see Figure 2). [11]

The RNNs or even better LSTMs architectures used for probabilistic language modeling can be reused in a more complex domain called sequence to sequence modeling for Neural Machine Translation (NMT) from one language to another. Here one first tries to learn a fixed-dimensional input representation from an input sequence using an encoder architecture based on an LSTM. The so-called context vector is then decoded by a second LSTM into a new sequence of words preserving the grammar but having a different meaning. [27]

Sequence to sequence models are introduced together with the LSTM architecture in subsection 2.5. Here the connection to evolution theory can be drawn. Ribonucleic Acid (RNA) sequences made of a concatenation of nucleotides¹ can be represented textually using the FASTA format [8]. A Seq2Seq model can then transferable be applied in the domain of RNA

¹We restrict the representation of nucleotides solely to their nucleobases parts consisting

sequences to model how RNA-based viruses change their structure to avoid the detection by the human immune system but still to preserve their infectivity and evolutionary fitness [13].

2.2 GISAID EpiFlu Data Platform

The data needed for this project consists of the SARS-CoV-2 genome data and the corresponding metadata. The most complete database is hosted by the GISAID initiative [2]. It is a public-private partnership between the initiative itself and the governments of Germany (official host of the platform), Singapore and the United States of America. [1]

The main aim of GISAID is to enable the fast sharing of epidemic and pandemic virus data. In comparison to other databases such as GenBank the genome sequences are shared fast through GISAID. For enabling this fast sharing the GISAID database is only usable after registration and the redistribution of the data is forbidden, whereas the GenBank database is publicly available. [25]

Thus, until september 2021 over three million SARS-CoV-2 genome sequences are available through GISAID, whereas only about one million SARS-CoV-2 genome sequences are available through GenBank. Figure 3 shows the selection view of the GISAID database. [1, 17]

2.3 Domain-Specific Methodologies to create Evolutionary Datasets for Mutation Prediction

From databases such as GISAID or GenBank an unordered list of genome sequences and metadata can be downloaded. For training a model to detect changes and predict future mutations the dataset must contain parent-child pairs.

Berman et al. [6] have done this based on two steps. First, they created a phylogenetic tree from the genome sequences. Therefore first all genome sequences are aligned. The initial phylogenetic tree is calculated by Fasttree using the approximate maximum likelihood method. In the next step, the tree is refined by a generalized time-reversible (GTR) model and a gamma model. The final phylogenetic tree is achieved by optimizing the branch length and the used models for tree generation. Like in standard phylogenetic trees each leaf node corresponds to one data instance, whereas the inner nodes are inferred from the other data instances. Secondly based on the rooted

of the distinct nucleobases guanine, adenine, cytosine and thymine. We therefore do not include the phosphate group and the five-carbon sugar components.

The screenshot shows the GISAID database search interface. At the top, there are tabs for Registered Users, EpiFlu™, EpiCoV™ (which is selected), EpiRSV™, and My profile. Below the tabs, there are links for EpiCoV™, Search, Downloads, and Upload. The main area is titled 'Search' and contains several search filters: Accession ID, Virus name, Location, Host, Collection, Submission, Clade, Lineage, Substitutions, Variants, and a date range from 'to'. There are also checkboxes for 'complete', 'high coverage', 'low coverage excl.', 'w/Patient status', and 'collection date compl.'. Below the filters is a table with columns: Virus name, Passage date, Accession ID, Collection date, Submission date, Length, Host, Location, and Originating location. The table lists 12 entries of hCoV-19/Mexico/ZAC_IBT_IMSS_2588/202 variants. At the bottom of the table, it says 'Total: 3,364,537 viruses' and has navigation buttons (">< < | 1 | 2 | 3 | 4 | 5 | > >>). To the right are buttons for 'Select', 'Analysis', and 'Download'.

Figure 3: GISAID database [own screenshot]

phylogenetic tree, parent-child pairs are calculated. Therefore first a marginal ancestral sequence reconstruction is executed based on the phylogenetic tree from the previous step. Now also the inner nodes correspond to data instances. In the next step BioPython's [9] Bio.Phylo package is used to determine parent-child pairs. From each edge, one parent-child pair is generated. [6]

Mohamed et al. [18] used existing datasets from Ogali et al. [21] in their work about mutation prediction. Ogali et al. [21] performed a phylogenetic analysis. After selecting only genome sequences with high quality, they removed duplicate sequences and added reference sequences to the dataset. Now the alignment is executed. In the next step MEGA (Molecular Evolutionary Genetics Analysis) is used to create the phylogenetic tree. This is done by using the maximum likelihood method, the best-fit general time-reversible (GTR) model and gamma-distributed rate variation.

Hadfield et al. [12] have created the nextstrain project. It aims to examine pathogen's spread and evolution by providing a viral genome database, a bioinformatics pipeline for phylodynamics analysis and a visualization platform. Steps executed by the bioinformatics pipeline for phylodynamic analysis are: "subsampling, alignment, phylogenetic inference, temporal dating of ancestral nodes and discrete trait geographic reconstruction, including in-

ference of the most likely transmission events" [12]. For calculating the phylogenetic tree, TreeTime's maximum likelihood method is used. [12]

2.4 Previous Work on Mutation Prediction

Even before the rise of Covid-19, there had been studies trying to predict mutations of RNA viruses. In the collection of [29, 30, 31] the authors predict the mutation positions in hemagglutinins from influenza A virus using logistic regression and plain neural networks and then use the resulting amino acid mutating probabilities to derive possible mutated amino acids. The same approach is further used for H5N1 neuraminidase proteins.

Salama et al. [23] proved that nucleotides in an RNA sequence can change based on their local neighborhood. Neural networks are used to predict new strains of the Newcastle virus and subsequently, a rough-set theory based algorithm is introduced to extract the according point mutation patterns.

Mohamed et al. [18] used a more modern sequence to sequence approach based on LSTMs to learn nucleotide mutations between time-series species of H1N1 Influenza virus and the Newcastle virus as mutations can also be influenced by long-distance relations of amino acids. Therefore one-hot encoded RNA sequences of a parent generation preprocessed to words are given as an input and the output is the predicted offspring generation evaluated by accuracy to the compared true offspring generation. The achieved accuracy in this paper is questionably high with 98.9% on the H1N1 Influenza virus and 96.9% on the Newcastle virus, possibly because of overfitting to the few 4.609 samples for H1N1 Influenza virus and only 83 for the Newcastle virus. The approach in this paper therefore tries to increase the number of samples available for training when building the dataset.

The approach in this paper will neither use any of the just mentioned architectures but uses a transformer-based architecture coupled with a GAN-style training architecture. A short introduction into Seq2Seq models and the underlying LSTM components shall be given to better point out the architectural decisions made.

2.5 Sequence-to-Sequence Models based on Long Short-Term Memory

The original LSTM unit was introduced by Hochreiter and Schmidhuber in [14] and can be used for language modeling instead of using plain RNNs to prevent running into vanishing or exploding gradient problems [26]. The architecture of an LSTM is shown in Figure 4.

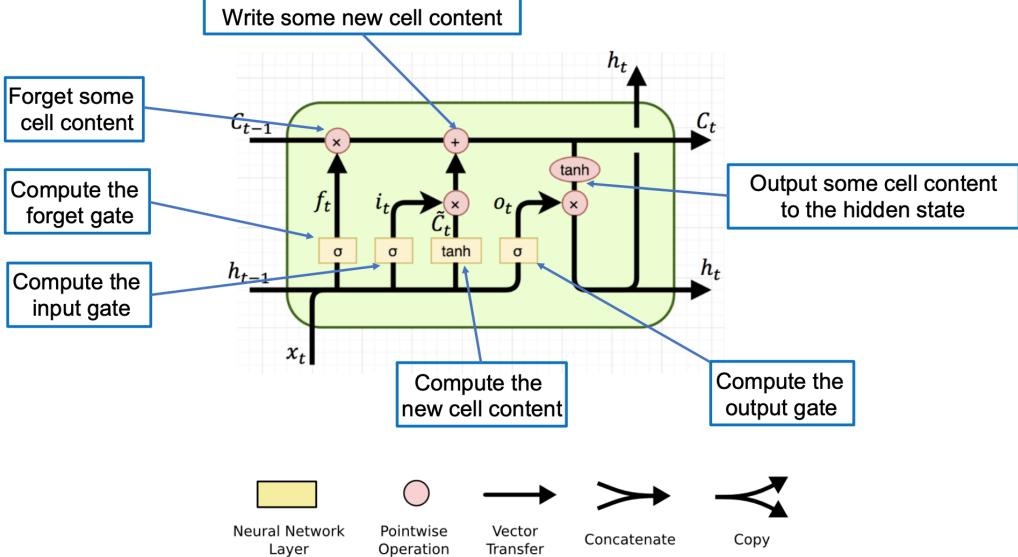


Figure 4: Architecture of an LSTM [11]

It consists of a hidden state h_t and an additional cell state c_t . The cell state stores long-term information and is used to derive a new hidden state. Information flows through three different gates inside the LSTM. The forget gate is used to control which parts of the cell state are potentially carried on to the next time step, the input gate is responsible to decide which parts of the cell state should be updated and the output gate determines what is being passed on as the new hidden state. All three gates depend on the previous hidden state and the current input. They provide factors limited to the interval $[0,1]$ by the sigmoid function and are multiplied with the cell state, the changes to be added to the cell state and the new hidden state derived from the cell state. Through the cell state, an LSTM therefore makes it possible to capture long-distance dependencies. [11]

Sutskever et al. [27] introduced Seq2Seq learning following a multi-layer encoder-decoder style model architecture. One layer consists of one LSTM that is used as an encoder to learn a large fixed-dimensional vector representation of a size-unrestricted input sequence called the context vector. This vector consists of the last cell and hidden state of the encoder and incorporates the structure of the input sequence helping the following decoder LSTM to provide qualitative predictions for the output sequence. The second LSTM therefore serves as a beam search² decoder to map the context vector

²Do not choose the most probable word but the B most likely word hypothesis and pass them to the next time step in the LSTM. Whichever hypothesis results in the lowest loss is kept. To avoid combinatorial explosion limit the beam depth size.

to a corresponding output sequence whose length does not need to match with the length of the input sequence. The output probability distribution is therefore given by the equation

$$p(y_1, \dots, y_{T'} | x_1, \dots, x_T) = \prod_{t=1}^{T'} p(y_t | v, y_1, \dots, y_{t-1}) \quad (4)$$

with v being the context vector. Using an LSTM is preferred over a normal RNN as it is used to capture the long-range temporal dependencies of the input data. The encoder-decoder architecture uses four layers in total partitioned onto four Graphics Processing Units (GPUs). A corpus of 160k words for the input sequence and another one of 80k words for the target sequence was used to create the word embeddings of dimension 1000. Unknown words were replaced by an *UNK* token. The Seq2Seq model approach was evaluated for NMT and reached a 34.81 Bilingual Evaluation Understudy (BLEU) score. One finding during training was that reversing the input sequence introduces many short-term dependencies as the minimal time lag of the problem is reduced making optimization easier. [27]

2.6 Applying Generative Adversarial Networks

Using a plain Seq2Seq model for mutation prediction does not necessarily guarantee that the generated sequences are evolutionary offsprings of a parent generation as not being included as is in the ground truth data. The generated sequences might occur realistic and biologically relevant, but a plain sequence to sequence architecture does not inherently check for natural parental descent and therefore does not make sure whether the predicted mutations lead to improved fitness. Berman et al. [6] developed a novel Seq2Seq framework based on the GAN idea to predict genetic mutations and future biological populations of the influenza virus (see Figure 5). MutaGAN describes a Seq2Seq generator within an adversarial framework that predicts protein sequences augmented with possible mutations. By using a Seq2Seq generator and a discriminator specialized on separating fake evolutionary mutations from real ones, one can then guarantee to a certain degree that the evolutionary parent-childhood coherence is given. In MutaGAN a mutation is considered correct if the change in amino acid and location within the RNA sequence is equal to the parent’s true offspring. [6]

MutaGAN’s generator consists of a Seq2Seq model. The encoder is built of a bidirectional LSTM and the resulting context vector of dimensionality 512 is combined with noise from a standard normal distribution. The decoder LSTM predicts the resulting sequence of amino acids greedily using the *argmax()* of the *softmax()* output for every time step. The discriminator

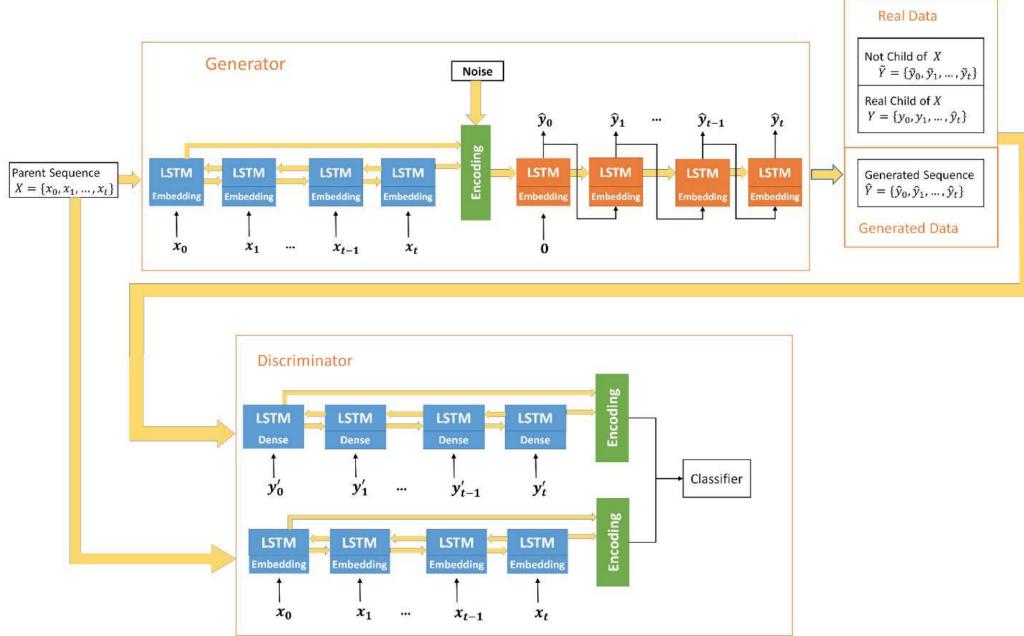


Figure 5: Architecture of MutaGAN [6]

is trained on three different parent-child pair configurations to optimally compete against the generator, real parent-child pairs, real parent and generated child pairs and pairs of real sequences that are not parent-child pairs. Two bidirectional LSTM encoders sharing the same weight matrix as the encoder of the generator produce a fixed-dimensional encoding of the provided parent-child pair used for classification of evolutionary descent by a plain neural network. The child encoder in the discriminator uses a dense layer having the same weight matrix as the embedding layer of the encoder of the generator. The dense layer is required to directly input the predicted child sequence as its probability distribution rather than the final output after applying the *argmax()* as it would not enable backpropagation to train the generator. When a true child sequence is given a simple one-hot encoding is used. [6]

Interestingly Berman et al. [6] state planned improvements by utilizing bigger datasets acquired from the GISAID database EpiFlu and by using more length-robust attention-based models to directly work on nucleotide sequences instead of sequences of amino acids. Therefore transformers and the attention mechanism are introduced in the following section.

2.7 Transformer and Attention Mechanism

Using a Seq2Seq model as introduced based on an encoder-decoder architecture of LSTM cells has some drawbacks. First, the input needs to be processed sequentially in time steps which makes parallelization difficult and increases training time, especially for longer sequences. Furthermore, the hidden and cell state vector passed through every timestep tries to encode information of *all* previous time steps without knowing which information is especially important for the current time step. This not only makes long-distance dependencies hard to capture but also might not get the prioritization of the previous time step inputs right. [5, 28]

To tackle these problems the so-called transformer architecture was introduced in [28]. It feeds an entire sequence into the encoder to be processed in parallel denying any concept of recurrence or convolution. Only using a so-called self-attention mechanism the transformer makes sure that during processing every input position, receives the information that is most important to it. This way modeling long dependencies becomes much easier compared to LSTMs as the view on the input sequence is more global. In this architecture, the encoder also passes all computed hidden states of every position to the decoder, which therefore can generate the target sequences based on more semantically enriched features and also in parallel for every position. [28]

The transformer architecture achieved state-of-the-art quality results with a BLEU score of 41.8 on the WMT 2014 English-to-French translation task (cf. BLEU score of [27] was 34.8), while still being much faster to train due to its parallelization capabilities. State-of-the-art results are already achieved after just twelve hours of training on eight P100 GPUs. As this is still far beyond the scope and resources given for this project, this project provides a proof-of-concept transformer model trained for far shorter and fewer sequences as one would need to predict entirely new RNA sequences. [28]

First the transformer architecture should be introduced in the following:

The transformer consists of an encoder and a decoder part just as normal Seq2Seq models. One layer of the encoder contains a self-attention layer before passing the hidden state representations of an input sequence to a feed-forward neural network consisting of two linear transformations separated by Rectified Linear Unit (ReLU). Thus it makes sure that semantic and contextual dependencies between different positions in the input sequence are also modeled. Also, skip connections are added for both components with subsequent layer normalization. [4, 28]

To calculate the self-attention output of a sequence, a query Q , key K

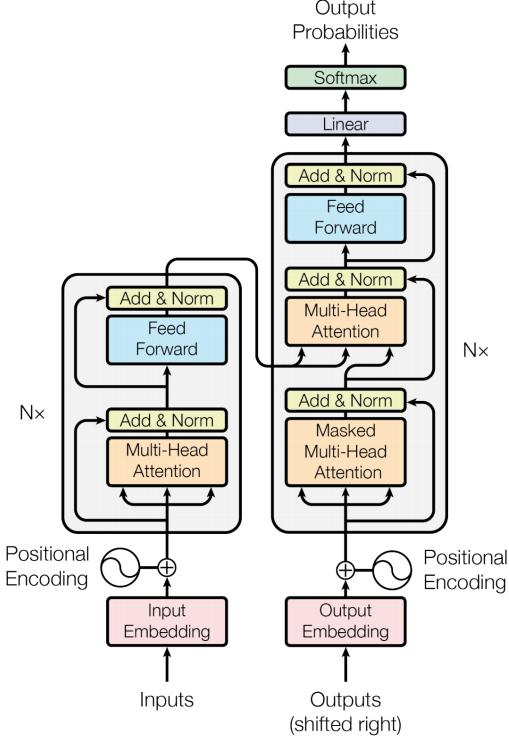


Figure 6: Architecture of the transformer [28]

and value V matrix is calculated through the multiplication of the sequence representation³ with learned transformation matrices. Each row of the three received matrices corresponds to a specific input position. The output of the self-attention layer is then calculated through

$$(ScaledDot - Product)Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V. \quad (5)$$

The dot product between Q and K calculates scores that are used as weighting factors for V . This describes for every input position how relevant all other input positions are for each of the hidden state encodings. The scores are divided by the square root of the dimensionality of the query/key values of the hidden state representations, which is chosen to be 64 (square root eight), to receive more stable non-vanishing gradients. The softmax guarantees that the positional scores lie between zero and one and add up to one. Note that the score of the current input position itself will most likely have the highest score. [4, 28]

³A matrix containing the hidden state representations of every input position

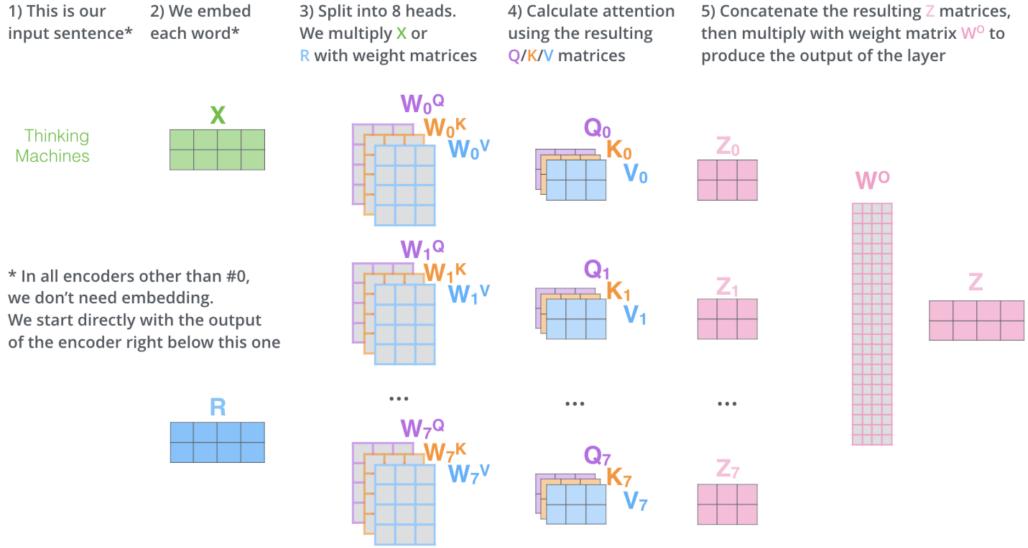


Figure 7: Multi-head attention architecture [4]

One can expand self-attention towards multi-head attention. In this configuration, each head produces its own query, key and value matrices and therefore its own self-attention output in its own representation subspace. The resulting hidden state is calculated by concatenating each hidden state output of the different heads and multiplying it with a contraction matrix that outputs a hidden state vector in its single-head dimensionality. Multi-head attention is needed to be able to better focus on multiple parts of the input sequence when encoding an input position, otherwise, it might happen that most of the focus is set on the input position itself. Vaswani et al. [28] used 8 heads. [4, 28]

Overall six layers of encoder and decoder components having individual weight matrices are stacked on top of each other to capture low-level as well as high-level features. The dimensionality of the hidden state is 512. Note that each position of the input sequence is encoded on an individual path through the transformer, which makes parallelization possible compared to LSTM based architectures. [4, 28]

The decoder uses almost the same architecture as the encoder but contains an additional multi-head attention component that integrates all the hidden state encodings from the encoder. Therefore the final hidden state encodings resulting from the top most encoder are transformed into key and value attention vectors and given to the integrating multi-head attention component of every decoder. Once again this helps to better focus on the relevant input positions for the current output position based on hidden state information

extracted from the encoder. Note that the decoder works in a sequential manner again meaning that the input to the decoder are the already decoded sequence tokens of previous time steps, future input positions are masked away. The first multi-head attention component therefore captures inter-dependencies in the generated output sequence, the second one, the so-called encoder-decoder attention, enriches the hidden state encoding with the information given from the encoder hidden states. Finally, a linear layer mapping the hidden state vector to a logits vector with the dimensionality of the given output dictionary and a *softmax()* function produce the output sequence using the greedy or beam search approach. The *EOS* token symbolizes the end of the output to be generated. [4, 28]

Input to the transformer are word embeddings also of fixed dimensionality of 512. Onto each word embedding a positional encoding following a specific pattern is added. These patterns are either learned or fixed and make sure that the distances of the input positions are projected onto the queue, key and value representations to be utilized during scaled dot-product attention. This also enables self-attention to be scaled to unseen lengths of sequences. Transformers usually contain an upper bound of sequence length to be processed, which is different from recurrent architectures that can process inputs of arbitrary length. Input sequences are usually padded up until the specified sequence length, the used *PAD* tokens are masked to be excluded from the self-attention components. [4, 28]

3 Approach

Figure 8 visualizes the steps of our approach, which are described in detail in the following chapters.

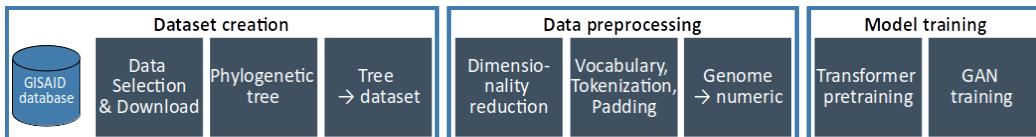


Figure 8: Approach pipeline [**own representation**]

3.1 Dataset Creation

As described in chapter 2.3 a dataset consisting of parent-child genome pairs is needed to learn a machine learning model for mutation prediction. The steps to create this dataset are described in the following chapters.

3.1.1 Raw Data Selection from GISAID

In the first step of the dataset creation, the raw genomes and their metadata are downloaded from GISAID. Therefore, a selection of a suitable subpart of the data is necessary. The GISAID platform only allows the download of 5000 genomes at once. That is why the decision was made to focus the analysis on Germany (for one country the selection of < 5000 genomes is rather simple compared to the selection of data from multiple countries). By looking at the latest "Report on virus variants of SARS-CoV-2 in Germany" from the Robert Koch Institut (RKI) a time period was chosen. Figure 9 shows the distribution of different variants over time. Starting from calendar week 18 the Delta variant gradually displaces the previously widespread Alpha variant.

That is why the data from 04.05.2021 (calendar week 18) until 15.07.2021 (calendar week 28) was selected for this study. This results in a raw dataset of 35.818 genomes. Each genome consists of a FASTA file containing the genomes sequence and a record in a .tsv file with the metadata.

3.1.2 Generation of a Phylogenetic Tree

As described above parent-child genome sequence pairs are needed to be able to train the model. Therefore one needs to evaluate the ancestral relationships between the genome sequences. As described in chapter A.1.4

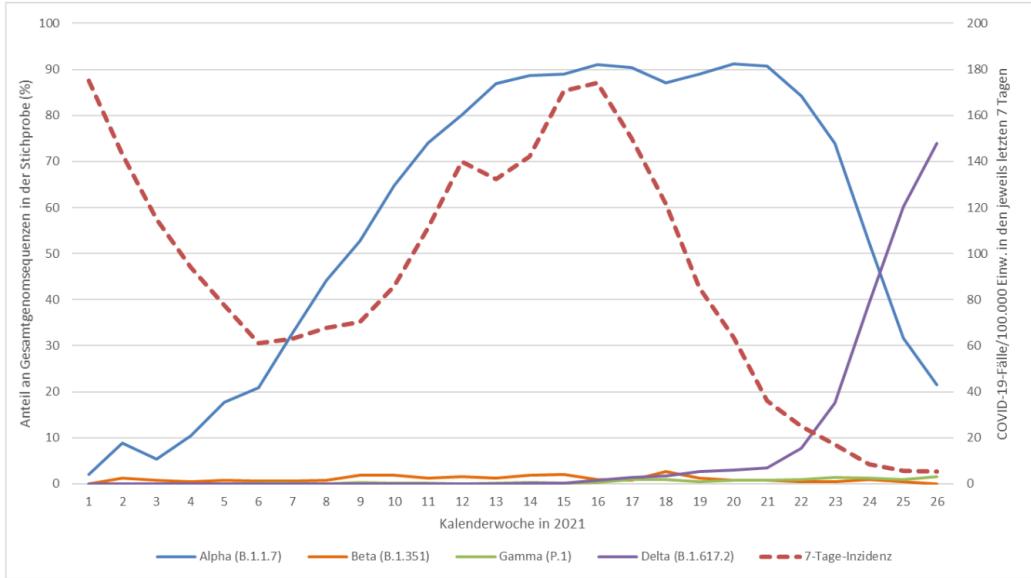


Figure 9: Variant distribution in Germany over time [3]

this is done in biology through phylogenetic trees. For the data preparation (e.g. alignment) and the calculation of the phylogenetic tree the nextstrain pipeline is used [12]. The pipeline was executed on a local computer with Intel Core i7-7500U (4*2.70GHz), NVIDIA GeForce 940MX and 16 GB RAM. Unfortunately, after two days of calculation, the nextstrain pipeline exited with an out-of-memory exception. That is why the amount of data was decreased. This was done through iterative reduction of the dataset size. The largest locally computable dataset consists of 11.773 genome sequences.

The generated phylogenetic tree can be seen in figure 10.

3.1.3 Phylogenetic Tree to Dataset

Based on the calculated phylogenetic tree the final dataset is generated. Therefore the patristic distances between all leaf nodes are calculated. From the resulting patristic distance matrix for each leaf node, the most related next leaf node can be evaluated. The two leaf nodes with the smallest distance are the closest relatives to each other and therefore represent a parent-child pair. After the evaluation of which leaf nodes are one parent-child genome pair, the question of which node is the parent and which node is the child is clarified. Therefore the two nodes are sorted by their date, which is stored in the metadata file. The older genome sequence is the parent and the younger genome sequence is the child.

For calculating the patristic distance matrix Dendropy (version: 4.5.2)

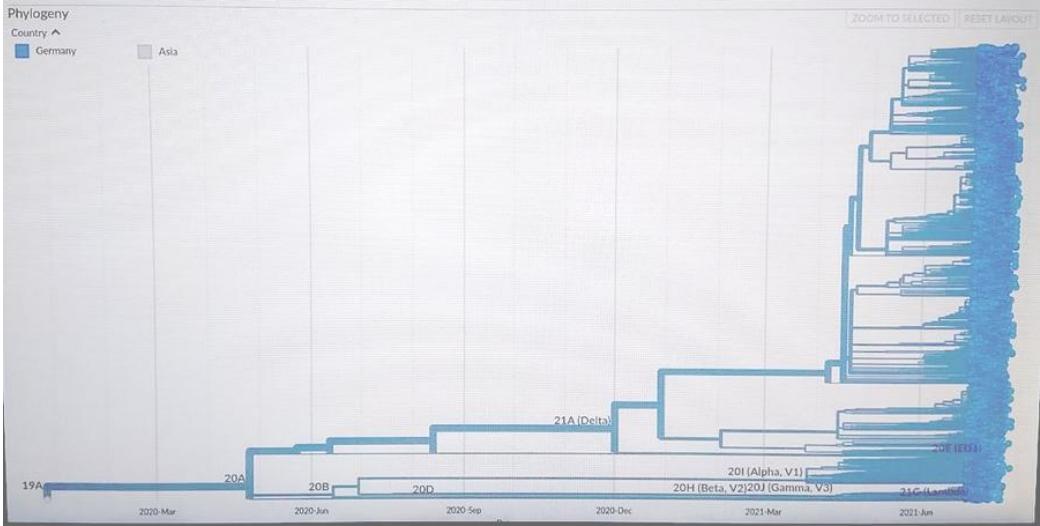


Figure 10: Phylogenetic tree of our dataset [**own representation**]

[10] was first used. Again this led to problems due to limited RAM on the local computer. As a consequence, a switch to PhyloDM (version: 1.3.1) [19] was done, which is a library for calculating patristic distances in a memory and time-efficient manner. Finally, a dataset containing parent-child genome sequences was received.

3.2 Data Preprocessing

3.2.1 Dimensionality Reduction by selecting Subpart of the Genome

The full SARS-CoV-2 genome consists of 29.904 nucleotides. Due to the limited computing capacity on the local computers, the dimensionality of the dataset had to be reduced. This is done by selecting a subpart strang based on the gained data insights from chapter 4.1. In figure 12 parts of the genome with lots of mutations are visible. Based upon this 99 nucleotides (33 codons) in the range between position 21.800 and 21.899 were selected. For simplicity, the selected strang sequences are treated as if they were full genome sequences in the following sections.

3.2.2 Transformation of Genome Sequence to Numeric Model Input

Input to the model are numericalized codon sequences with each codon being composed of three nucleotides (one amino acid). Therefore one needs to build a vocabulary with all codons existing in the dataset that assigns each codon a unique number. Each genome string is therefore tokenized to codons and so

far unseen codons are added to the vocabulary. The vocabulary also contains special tokens:

- $\langle SOS \rangle$: start of sentence
- $\langle EOS \rangle$: end of sentence
- $\langle UNK \rangle$: unknown
- $\langle PAD \rangle$: padding

A codon is encoded as unknown in the case at least one nucleotide is unknown or at least one padding character is contained. Only in the special case of three padding characters, the padding token is returned.

To integrate the dataset into the training routine, PyTorchs *DataLoader* is used (see subsection 3.4). Therefore a *CustomGISAIDDataset* class is derived from PyTorchs *Dataset* class that can be utilized by the *DataLoader*. It can be configured to return the training or test dataset instances. Each instance is already numericalized by vocabulary index lookups for each codon, padded with the numericalized $\langle SOS \rangle$ and $\langle EOS \rangle$ token and transformed to a PyTorch tensor.

3.3 Model Architecture

The connection of modeling evolution theory of virus mutations to probabilistic language modeling has already been explained in subsection 2.1. To create such a probabilistic language model an approach closely related to the one introduced by Berman et al. [6] as explained in subsection 2.6 is chosen with the novelty that attention-based transformers are used instead of LSTMs⁴. To implement the GAN framework, PyTorch version 1.9.0 is used as it provides full-featured classes for transformers, embedding layers and other architectural components needed compared to other frameworks like Keras. This way one can specialize on the training routine instead of implementing a custom transformer architecture, which might be too error-prone⁵.

The architecture of the generator was shown in Figure 6. The numericalized codon sequences of the source parent and the target child genome sequence first need to be transformed to input embeddings of size 32 for each codon. This size was chosen experimentally and seemed to work best

⁴This improvement was also named in the conclusion of the MutaGAN paper [6].

⁵In case one is interested in how to implement a custom transformer architecture, see <http://nlp.seas.harvard.edu/2018/04/03/attention.html>.

in practice as it only needs to model a vocabulary of size 40 for the parent codon vocabulary and 40 for the child codon vocabulary. To those input embeddings, positional encodings are added that are formed by embedding each index in a sequence of positional indices for each codon in a genome instance. To both input sequences dropout layers are applied with a dropout rate of 10%. The transformer then predicts the child genome sequence as a hidden representation of dimensionality 64 for each codon. The hidden representation is the output of the feed-forward layer in a transformer block, its dimensionality was also chosen according to best experimental performance. In total three encoder and decoder layers are used each having eight heads. The transformer also uses a dropout rate of 10%. The final feed-forward network reduces the dimensionality from 64 back to the size of the vocabulary for each codon and applies the softmax function to receive the final output probabilities. Using greedy decoding the codon with the highest probability is returned in the evaluation phase as well as during training for simplicity reasons (instead of using a beam search approach). Note that padding tokens or so far unpredicted child tokens are masked away to not be considered when calculating the loss during the training routine.

The architecture of the discriminator is kept similar to Figure 5 using PyTorch’s *TransformerEncoders* instead of LSTMs. Again one transforms the numericalized parent and either true or generated child genome sequence to embeddings of size 32 for each codon and adds the positional embeddings before applying dropout. Instead of using an embedding layer for the true or generated child genome sequence, one uses a dense layer instead, to directly feed in the probability distribution of generated child sequences. The true child sequence therefore needs to be one-hot encoded before being provided to the discriminator to match the shape. The transformer encoder then produces an encoding of dimensionality 64 for each codon. During the training phase dropout is once again used with a rate of 10%. Note that until this step the same weights are used as in the generator. In the following step, the two hidden representations of the parent and true or generated child genome sequence are concatenated along the embedding dimension. The subsequent feed-forward network uses three blocks of dropout, batch normalization and a linear layer in this order with ReLU activation. The ReLU layer in the last block is replaced with a sigmoid layer that makes sure each codon receives a score between zero and one that states how likely its a true or generated codon. Once again padding tokens are masked away during the training phase.

3.4 Training Process

The training phase is two-fold and is performed on a single GeForce GTX 1080 GPU. First, a pretraining phase for the generator needs to be performed in order to produce the first reasonable child genome sequences. Also, the weights of the embedding layers and the encoder part of the generator will be reused for the discriminator in the second training phase. For the pretraining phase, the cross-entropy loss function typical for language modeling is used, the learning rate is set to 0.005. Teacher forcing is used in the initial pretraining phase to accelerate the training process by leading the generator to the right direction towards reasonable child genome sequences. In the main training phase, the generator and discriminator compete against each other. When training the discriminator, the generator first generates a child genome sequence given a parent genome sequence. The true and generated child genome sequence is then given to the discriminator. Here the binary cross-entropy loss is used to model the two parts of the classical GAN loss function:

$$\underset{G}{\operatorname{argmin}} \underset{D}{\operatorname{argmax}} E_{x \sim p^*(x)}(\log(D(x))) + E_{z \sim p(z)}(\log(1 - D(G(z)))) \quad (6)$$

PyTorch’s *BCELoss* function is given as:

$$l = -w(y \log(x)) + (1 - y) \log(1 - x) \quad (7)$$

We therefore pass $y = 1$ when calculating the left part of the loss using the true child genome sequence and $y = 0$ when calculating the right part of the loss using the generated child genome sequence. We use the same mechanism to calculate the loss for the generator, but maximize $\log(D(G(z)))$ instead of minimizing $\log(1 - D(G(z)))$. After experiencing mode collapse problems Berman et al. [6] switched to the Wasserstein loss function instead of the binary cross-entropy function and therefore also changed the last layer of the discriminator to be a linear activation function. As no mode collapse was discovered and PyTorch does not offer a Wasserstein loss function the binary cross-entropy loss function was kept. The learning rate for the generator was set to 0.0001 whereas the learning rate for the discriminator was set far lower to 0.00003 to avoid mode collapse. Besides providing true child genome sequences and generated child genome sequences to the discriminator, Berman et al. [6] also provide unrelated, but true child genome sequences to the discriminator. This way it is hoped to not only learn which mutations are relevant as inherently given by the training dataset but also to strengthen the knowledge of which mutations are especially related to a specific parent

genome sequence and which not. For time reasons the training dataset was not enriched with such true but unrelated child genome sequences.

Both training phases use batch gradient descent as their optimization method with a batch size of 32. Gradient descent is performed alternatingly between the generator and discriminator on every epoch. Gradients are clipped to a maximal norm of one to avoid exploding gradients. Furthermore, for both training phases, the Adam optimizer is used. A learning rate scheduler makes sure to reduce the learning rate by a factor of ten in case the loss has not improved for ten epochs. On every batch gradient descent, a tensorboard is updated to draw the loss curve. Every ten epochs a model checkpoint is saved. Pretraining is done for only 50 epochs whereas the main training phase runs for 300 epochs.

4 Experimental Results

4.1 Dataset Insights

4.1.1 Complete Dataset

The dataset consists of 9199 parent-child genome pairs. The distribution of how the parent and child sequences differ can be seen in figure 11. The majority of parent-child genome pairs differ in less than 200 characters. The number of completely equal parent-child genome pairs is 396.

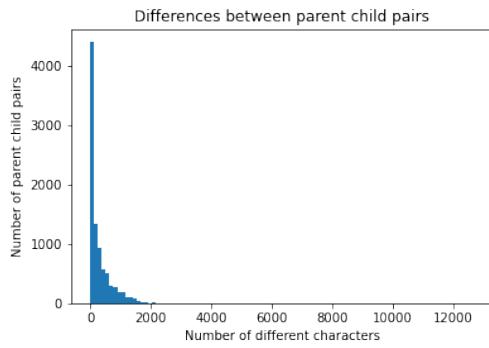


Figure 11: Distribution of the differences between parent and child sequences [own representation]

Figure 12 visualizes the distribution of mutations over the full SARS-CoV-2 genome.

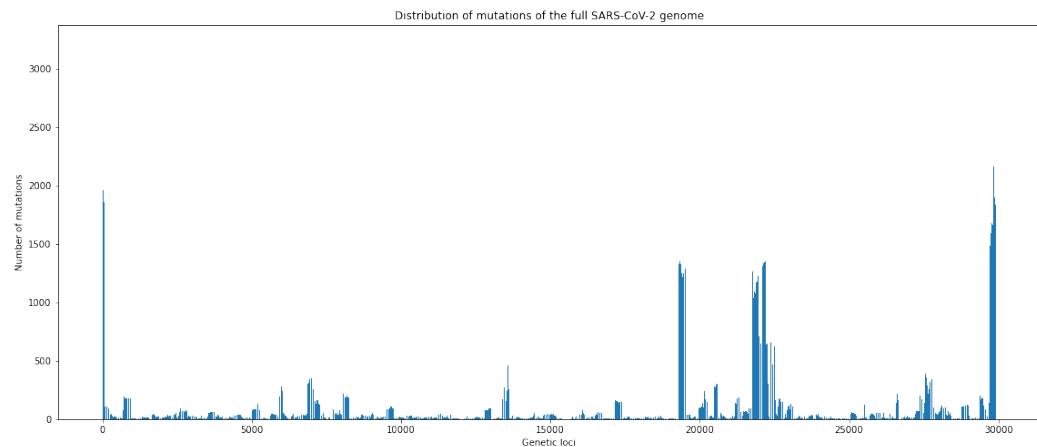


Figure 12: Mutated genetic loci [own representation]

4.1.2 Preprocessed Dataset

The complete dataset is preprocessed as described in chapter 3.2. The distribution of how the parent and child amino acids differ in the selected subpart can be seen in figure 13. The majority (6946) of parent-child genome pairs are equal. This is not optimal for training the model and is further discussed in chapter 4.2.

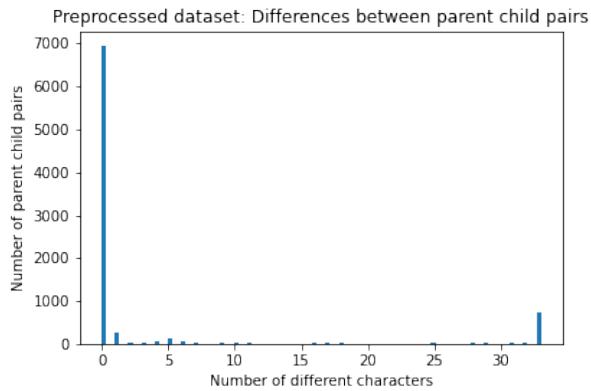


Figure 13: Preprocessed dataset: Distribution of the differences between parent and child sequences [own representation]

Figure 14 visualizes the differing amino acids and their positions of the selected subpart of the SARS-CoV-2 genome. They are distributed equally over the selected subpart.

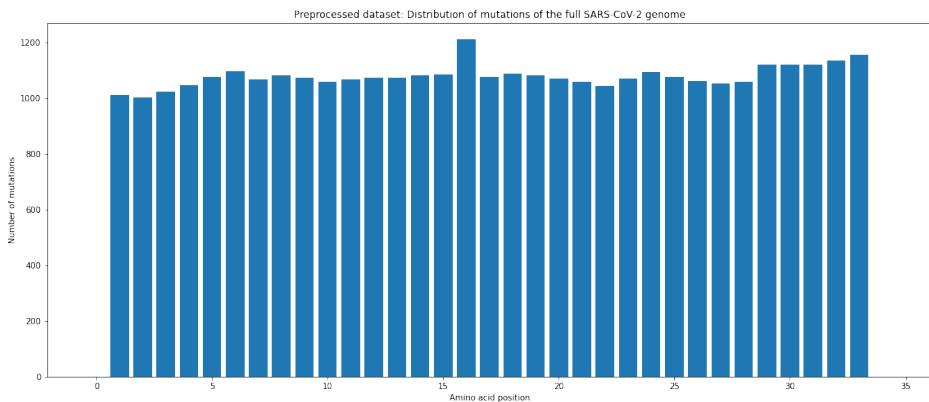


Figure 14: Preprocessed dataset: Differing amino acid positions [own representation]

4.2 Evaluation

4.2.1 Evaluation Criteria

To measure the success of the transformer model different evaluation criteria need to be introduced that capture how well the mutations observed in the data can be reproduced by the generator. A mutation prediction is considered correct in the case the prediction matches in codons and locations. Further to distinguish are partially correct mutation predictions in case mutations are observed at the right locations but in different codons. Mutation predictions can further be considered as wrong in case they were not observed in the true child sequences or missed in case they could not be observed in the generated sequences.

A common metric used in the domain of NMT is the so-called BLEU score [22]. Its so-called modified n-gram precision captures common n-grams between the generated sentence and the true reference translation sentence independently of their position. Also, a sentence brevity penalty is incorporated into the score. The BLEU score ranges between zero and one with one being a perfect match to the reference translation. [22]

Using solely the BLEU score would mean less focus on the concrete mutation prediction at critical genome positions but rather a correct overall child genome sequence prediction. As the child genome sequence is mostly similar to the parent genome sequence, predicting a nearly identical child genome sequence is not too difficult and the BLEU score would already resolve to a high value at the early stages of training.

Therefore the Berman et al. [6] introduced several other metrics besides the BLEU score. For instance, the distribution of specific codon mutations between the ground truth data and the generated data was calculated to evaluate the similarity of both mutation profiles. Furthermore, a so-called sequence true positive rate is determined by capturing all observable predicted mutations for each parent genome sequence and comparing them to the list of actual existing mutations for each parent genome sequence in the ground truth data.

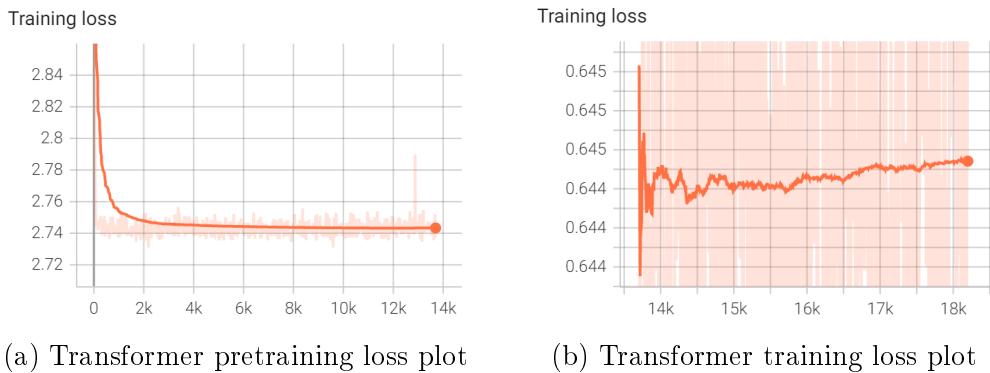
Besides the BLEU score, the sequence true positive rate is also utilized to evaluate the model using Google’s *Diff Match and Patch* libraries to identify the mutations⁶. The Levenshtein distance to determine the degree of similarity between a parent and generated child sequence is not used as it is not only partly covered by the BLEU score’s modified n-gram precision but also as at most two mutations appear for a given parent genome during

⁶For the Diff Match and Patch libraries, see <https://github.com/google/diff-match-patch>

the evaluation phase making the Levenshtein distance neglectable as it will most likely resolve in a very small value.

4.2.2 Training Phase and Evaluation

Figure 15a shows the loss plot for the transformer pretraining. The model converges rather fast. Similar is the actual GAN training (see 15b).



Whereas the trained model achieves a relatively high BLEU score with 87.42%, the sequence true positive rate only amounts to 31%. Thus one concludes that the model generates reasonable and close genomes from given parent genome sequences, but at first glance develops only partly the capability to predict the actual and correct mutations given in the ground truth data. When further analyzing the rather low sequence true positive rate, one recognizes that a lot of false predictions were made that can be traced back to cases where the parent genome sequence consists of a series of $<UNK>$ tokens that are replaced by actual codons from the transformer. As these kinds of mutations are not present in the ground truth data, whose matching child pairs still contain the $<UNK>$ series, one can think of an additional capability of the trained transformer model as a gap filler of unknown genome sequence parts of parent genomes, which is even beneficial. Other mutations were captured really well, e.g. from codon *TTG* to *CTG* at strand location 45 or from codon *TGC* to *TAC* at strand location 13. During evaluation, at most two mutations on the same parent in one prediction appeared, matching to the actual ground truth. Some mutations in the ground truth data were also missed, especially in case they were not too frequent. Only 25 unique parent sequences were considered in the test dataset further reducing a diverse mutation prediction capability of the model. But as data selection and preprocessing is very time and especially resource extensive⁷, the scope

⁷Phylogenetic tree reconstruction might take several days until completion.

of the project and its data used is intentionally kept minimal. In this sense, one can conclude that the model first not only generates reasonable genome offsprings but also features predicting noticeable codon mutations observable the ground truth data. As a comparison, MutaGAN achieves a higher BLEU score of 97.46% but an unweighted sequence true positive rate of just 21% [6]. When up-scaling the number of genome sequences to use and doing more strict filtering of genome sequences to increase the data quality, predicting virus mutations using transformers is very well feasible.

5 Conclusion

Coming to a conclusion, the main research question was whether a machine learning model can be trained to predict the next possible SARS-CoV-2 mutations. The short answer to this is, that the trained model is definitely able to generate possible next genome sequences demonstrating a BLEU score of 87.42% and a sequence true positive rate of 31%.

Based on related works introduced in chapter 2, a model architecture consisting of a transformer-based GAN framework was chosen. This architecture was influenced by Berman et al. [6] and builds upon its improvement proposal to incorporate transformers instead of LSTMs as a central novelty in this paper leveraging the attention mechanism and parallel processing power of the transformer architecture to achieve higher quality predictions.

For answering the research question a new dataset based on raw data from GISAID was created. As part of this, a reusable pipeline for creating evolutionary datasets was created.

Another novelty is the application and training of such a machine learning architecture for SARS-CoV-2.

As an outlook, there are some improvements and further investigations we would like to evaluate in future work. As realized while generating data insights of the dataset, lots of our data instances are very similar to each other. To really model the worldwide development of SARS-CoV-2 mutation events, the raw dataset should be based on a subsampling of all available global data instances. Further improvements could be achieved by using a computing cluster for generating a larger dataset and for training the model on the full SARS-CoV-2 genome and not on a subset of the genome. Moreover, beam search instead of greedy decoding during the training phase could be further examined.

A Appendix

A.1 Appendix A: Basics from Molecular Biology

A.1.1 Genome Sequences

The human genome is encoded as DNA. The DNA carries genetic instructions on how DNA based organisms develop and behave. It is structured as a double helix and consists of two nucleotide pair combinations: adenine (A), thymine (T) and cytosine (C), guanine (G). [15, p. 8]

A.1.2 Protein Biosynthesis

The protein biosynthesis (also known as the central dogma of molecular biology) describes how proteins are produced based on DNA sequences. The proteins then produce characteristics (e.g. human hair color). So the DNA directly influences which human characteristics are developed. [24, p. 6]

The process of protein biosynthesis can be seen in Figure 17. The single steps are described in the following:

1. Transcription: The DNA is transcribed by the RNA-Polymerase into pre-messenger RNA. As part of this process, the nucleotide thymine is replaced by uracil (U). [24, p. 9]
2. mRNA Processing: The pre-messenger RNA is transformed into the messenger RNA by removing the non-coding regions (introns). So the main difference between DNA and RNA is that the DNA is structured as a double helix, whereas the RNA consists of a single strain. [24, p. 9]
3. Translation: Ribosomes translate the mRNA into amino acids. Three nucleotides (one codon) decode one amino acid. The matching which codons decode for which amino acid can be seen in Figure 18. [24, p. 9]

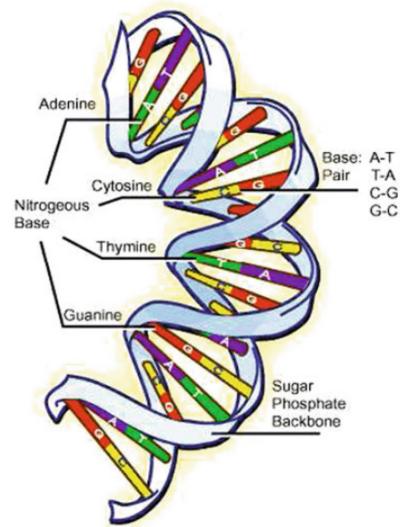


Figure 16: Deoxyribonucleic Acid (DNA) double helix [15, p. 8]

Central dogma of molecular biology

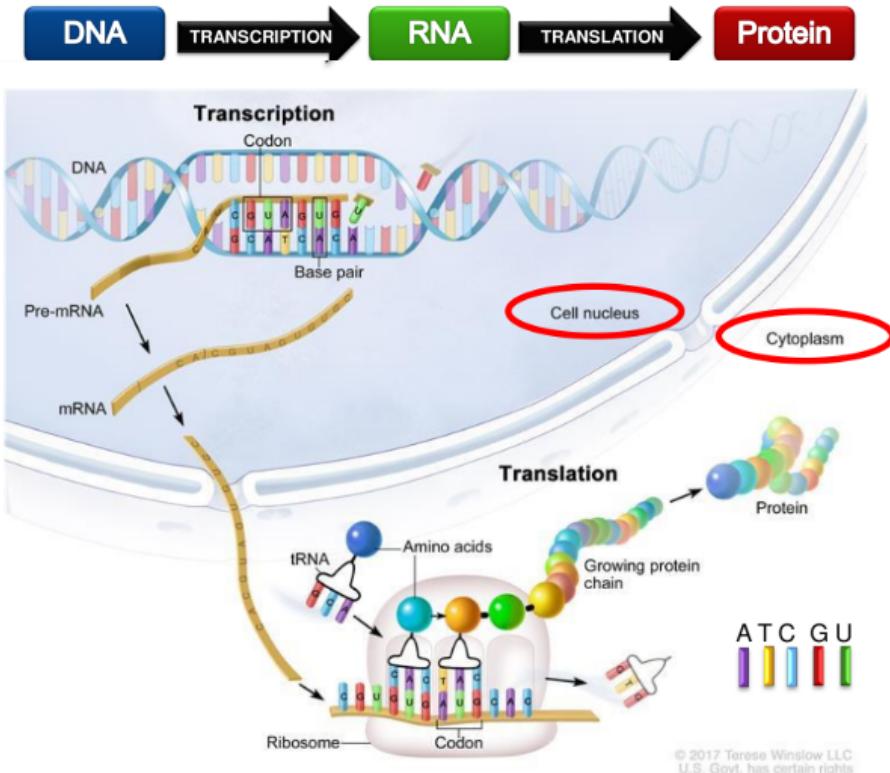


Figure 17: Protein biosynthesis [24, p. 6]

	U	C	A	G	
U	UUU Phe UUC UUA Leu UUG	UCU Ser UCC UCA UCG	UAU Tyr UAC UAA STOP UAG STOP	UGU Cys UGC UGA STOP UGG Trp	U C A G
C	CUU CUC Leu CUA CUG	CCU CCC CCA Pro CCG	CAU His CAC CAA Gln CAG	CGU CGC Arg CGA CGG	U C A G
A	AUU AUC Ile AUA AUG Met	ACU ACC ACA Thr ACG	AAU Asn AAC AAA Lys AAG	AGU Ser AGC AGA Arg AGG	U C A G
G	GUU GUC Val GUA GUG	GCU GCC GCA Ala GCG	GAU Asp GAC GAA Glu GAG	GGU GGC GGA Gly GGG	U C A G

Figure 18: Genetic code [24, p. 10]

A.1.3 Structure of SARS-CoV-2

RNA based viruses (like SARS-CoV-2) use RNA to encode their genetic information. That means that the genetic information is present as a single stranded RNA sequence. This sequence contains for example the information on which proteins form the surface structure of the SARS-CoV-2 virus and how this surface looks like. [20]

Based on Naqvi et al. [20] the four proteins forming SARS-CoV-2 are:

- spike (S)
- envelope (E)
- membrane (M)
- nucleocapsid (N)

A structural representation of SARS-CoV-2 and a host cell is visualized in Figure 19. [20]

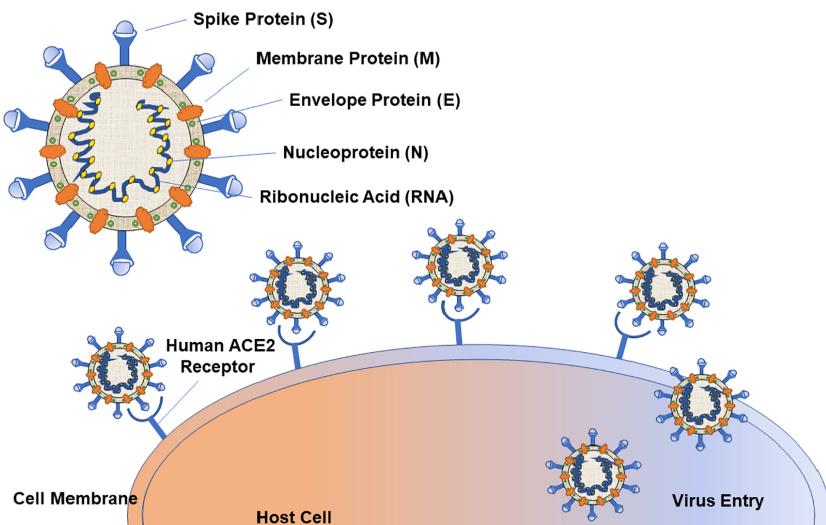


Figure 19: SARS-CoV-2 structure [20, p. 2]

Furthermore Figure 19 shows the RNA in the center. The SARS-CoV-2 RNA consists of about 30,000 nucleotides, which belong to different subgroups, as shown in Figure 20. Padded by a 5' Untranslated Region (UTR) at the beginning and a 3' UTR in the end (UTR are markers for the beginning and the end of protein-coding sequences), the middle part contains 12 functional Open Reading Frames (ORF), which contain the four coding regions for the proteins S, E, M and N.

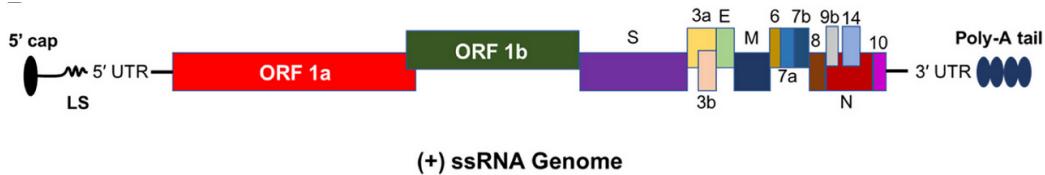


Figure 20: SARS-CoV-2 genome structure [20, p. 3]

A.1.4 Relationships between Biological Objects

One main question in biology is to determine relationships between objects, such as species or genome sequences. The observed objects are called taxa. Most commonly used for determining these relationships is the phylogenetic analysis which generates a phylogenetic tree. In this phylogenetic tree, each leaf node corresponds to exactly one taxon. The inner nodes of the tree are the inferred hypothetical ancestors, which are no taxa. The relatedness between different taxa can be evaluated by their distance. [7]

Figure 21 shows an example phylogenetic tree calculated based on the species genomes. On the right side, each leaf node (taxa) corresponds to one current species. From left to right one can see the inferred historical development from the ancestors to the current species. One can see, that Gorillas and Orangutans have earlier developed apart than the other species. Furthermore one can see, that the Chimpanzees and Bonobos share a most recent common ancestor (inner node next to both). That is why they are more related to each other. [16]

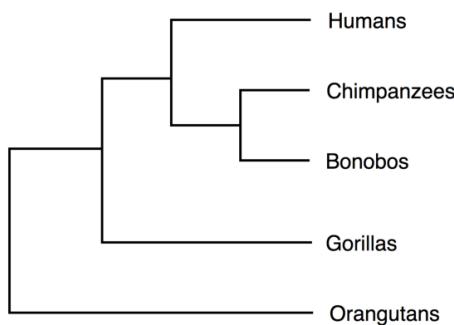


Figure 21: Example phylogenetic tree [16]

References

- [1] GISAID (editor). *GISAID*. Sept. 2021. URL: <https://www.gisaid.org/> (visited on 09/08/2021).
- [2] GISAID (editor). *GISAID - Mission*. GISAID - Mission. 2021. URL: <https://www.gisaid.org/about-us/mission/> (visited on 07/05/2021).
- [3] Robert Koch Institut (editor). “Bericht zu Virusvarianten von SARS-CoV-2 in Deutschland”. In: Bericht zu Virusvarianten von SARS-CoV-2 in Deutschland 18 (July 14, 2021), p. 19.
- [4] Jay Alammar. *The Illustrated Transformer*. 2018. URL: <https://jalammar.github.io/illustrated-transformer/> (visited on 08/18/2021).
- [5] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. “Neural Machine Translation by Jointly Learning to Align and Translate”. In: *arXiv* (2016). URL: <https://arxiv.org/abs/1409.0473> (visited on 08/19/2021).
- [6] Daniel S. Berman et al. “MutaGAN: A Seq2seq GAN Framework to Predict Mutations of Evolving Protein Populations”. In: *arXiv* (2020). URL: <https://arxiv.org/abs/2008.11790> (visited on 07/05/2021).
- [7] H.J. Böckenhauer and D. Bongartz. *Algorithmische Grundlagen Der Bioinformatik: Modelle, Methoden Und Komplexität*. XLeitfäden Der Informatik. Vieweg+Teubner Verlag, 2013. ISBN: 978-3-322-80043-5. URL: <https://books.google.de/books?id=2HL3BQAAQBAJ>.
- [8] P. J. A. Cock et al. “Biopython: Freely Available Python Tools for Computational Molecular Biology and Bioinformatics”. In: *Bioinformatics* 25.11 (June 1, 2009), pp. 1422–1423. ISSN: 1367-4803, 1460-2059. DOI: 10.1093/bioinformatics/btp163. URL: <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btp163> (visited on 09/20/2021).
- [9] Peter J. A. Cock et al. “Biopython: Freely Available Python Tools for Computational Molecular Biology and Bioinformatics”. In: *Bioinformatics* 25.11 (Mar. 2009), pp. 1422–1423. ISSN: 1367-4803. DOI: 10.1093/bioinformatics/btp163. eprint: <https://academic.oup.com/bioinformatics/article-pdf/25/11/1422/944180/btp163.pdf>. URL: <https://doi.org/10.1093/bioinformatics/btp163>.
- [10] *DendroPy Phylogenetic Computing Library — DendroPy 4.5.2 Documentation*. URL: <https://dendropy.org/> (visited on 09/17/2021).
- [11] Michael Gertz. “Text Analytics - Text Classification, Lecture Script”. University Heidelberg, 2020.

- [12] James Hadfield et al. “Nextstrain: Real-Time Tracking of Pathogen Evolution”. In: *Bioinformatics* 34.23 (May 2018), pp. 4121–4123. ISSN: 1367-4803. DOI: 10.1093/bioinformatics/bty407. eprint: <https://academic.oup.com/bioinformatics/article-pdf/34/23/4121/26676762/bty407.pdf>. URL: <https://doi.org/10.1093/bioinformatics/bty407>.
- [13] Brian Hie et al. “Learning the Language of Viral Evolution and Escape”. In: *Science* (2021). URL: <https://science.sciencemag.org/content/371/6526/284> (visited on 07/05/2021).
- [14] Sepp Hochreiter and Jürgen Schmidhuber. “Long Short-Term Memory”. In: *ResearchGate* (1997). URL: https://www.researchgate.net/publication/13853244_Long_Short-term_Memory (visited on 08/13/2021).
- [15] Nan M. Laird and Christoph Lange. *The Fundamentals of Modern Statistical Genetics*. 1st. Springer Publishing Company, Incorporated, 2010. ISBN: 978-1-4419-7337-5.
- [16] Vijini Mallawaarachchi. *Molecular Phylogenetics Using Bio.Phylo*. Medium. Mar. 24, 2018. URL: <https://towardsdatascience.com/molecular-phylogenetics-using-bio-phylo-57ce27492ee9> (visited on 09/09/2021).
- [17] National Library of Medicine (NCBI) (editor). *NCBI SARS-CoV-2 Resources*. URL: <https://www.ncbi.nlm.nih.gov/sars-cov-2/> (visited on 09/08/2021).
- [18] Takwa Mohamed et al. “Long Short-Term Memory Neural Networks for RNA Viruses Mutations Prediction”. In: *Hindawi* (2021). URL: <https://www.hindawi.com/journals/mpe/2021/9980347/> (visited on 08/11/2021).
- [19] Aaron Mussig. *Aaronmussig/PhyloDM: 1.3.1*. Version 1.3.1. Zenodo, Oct. 2020. DOI: 10.5281/zenodo.4089111. URL: <https://doi.org/10.5281/zenodo.4089111>.
- [20] Ahmad Abu Turab Naqvi et al. “Insights into SARS-CoV-2 Genome, Structure, Evolution, Pathogenesis and Therapies: Structural Genomics Approach”. In: *Biochimica et Biophysica Acta (BBA) - Molecular Basis of Disease* 1866.10 (2020), p. 165878. ISSN: 0925-4439. DOI: 10.1016/j.bbadi.2020.165878. URL: <https://www.sciencedirect.com/science/article/pii/S092544392030226X>.

- [21] Irene N. Ogali et al. “Molecular Characterization of Newcastle Disease Virus from Backyard Poultry Farms and Live Bird Markets in Kenya”. In: *International Journal of Microbiology* 2018 (Aug. 5, 2018). Ed. by Simona Nardoni, p. 2368597. ISSN: 1687-918X. DOI: 10.1155/2018/2368597. URL: <https://doi.org/10.1155/2018/2368597>.
- [22] Kishore Papineni et al. “BLEU: A Method for Automatic Evaluation of Machine Translation”. In: *ACM* (2002). URL: <https://dl.acm.org/doi/10.3115/1073083.1073135> (visited on 09/23/2021).
- [23] Mostafa Salama, Aboul Ella Hassanien, and Ahmad Mostafa. “The Prediction of Virus Mutation Using Neural Networks and Rough Set Techniques”. In: *EURASIP Journal on Bioinformatics and Systems Biology* (2016). URL: <https://bsb-eurasipjournals.springeropen.com/articles/10.1186/s13637-016-0042-0> (visited on 07/05/2021).
- [24] Dominique Scherer and Justo Lorenzo Bermejo. “Statistical Genetics and Genetic Epidemiology - SNPs and Other Common Variants”. lecture manuscript. lecture manuscript. University Heidelberg, 2021.
- [25] Yuelong Shu and John McCauley. “GISAID: Global Initiative on Sharing All Influenza Data – from Vision to Reality”. In: *Eurosurveillance* 22.13 (Mar. 30, 2017). ISSN: 1560-7917. DOI: 10.2807/1560-7917.ES.2017.22.13.30494. URL: <https://www.eurosurveillance.org/content/10.2807/1560-7917.ES.2017.22.13.30494> (visited on 09/08/2021).
- [26] Martin Sundermeyer, Ralf Schlüter, and Hermann Ney. “LSTM Neural Networks for Language Modeling”. In: *ISCA Archive* (2012). URL: https://www.isca-speech.org/archive/archive_papers/interspeech_2012/i12_0194.pdf (visited on 08/12/2021).
- [27] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. “Sequence to Sequence Learning with Neural Networks”. In: *arXiv* (2014). URL: <https://arxiv.org/abs/1409.3215> (visited on 05/15/2021).
- [28] Ashish Vaswani et al. “Attention Is All You Need”. In: *arXiv* (2017). URL: <https://arxiv.org/abs/1706.03762> (visited on 08/17/2021).
- [29] Guang Wu and Shaomin Yan. “Prediction of Mutations Engineered by Randomness in H5N1 Hemagglutinins of Influenza A Virus”. In: *SpringerLink* (2007). URL: <https://link.springer.com/article/10.1007/s00726-007-0602-4> (visited on 08/13/2021).
- [30] Guang Wu and Shaomin Yan. “Prediction of Mutations Engineered by Randomness in H5N1 Neuraminidases from Influenza A Virus”. In: *SpringerLink* (2007). URL: <https://link.springer.com/article/10.1007/s00726-007-0579-z> (visited on 08/13/2021).

- [31] Guang Wu and Shaomin Yan. “Prediction of Mutations in H1 Neuraminidases from North America Influenza A Virus Engineered by Internal Randomness”. In: (2008). URL: <https://link.springer.com/article/10.1007/s11030-008-9067-y> (visited on 08/13/2021).