

Heidelberg University  
Institute of Computer Science

Project report for the lecture Advanced Machine  
Learning

Prediction of the next SARS-CoV-2  
variants

<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

Felix Hausberger

# Contents

<b>0</b>	<b>Project Setup</b>	<b>2</b>
<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Fundamentals and Related Work</b>	<b>3</b>
2.1	From Language Models to modeling Evolution Theory . . . . .	3
2.2	GISAID EpiFlu . . . . .	3
2.3	Domain-Specific Methodologies to create Evolutionary Datasets for Mutation Prediction . . . . .	3
2.4	Domain-Specific Methodologies to create Evolutionary Datasets for Mutation Prediction . . . . .	3
2.5	Sequence2Sequence Models based on Long Short-Term Memory	3
2.6	Applying Generative Adversarial Networks . . . . .	3
2.7	Transformer and Attention Mechanism . . . . .	3
2.8	Other Techniques . . . . .	3
<b>3</b>	<b>Approach</b>	<b>4</b>
3.1	Dataset Creation . . . . .	4
3.2	Data Preprocessing . . . . .	4
3.3	Model Architecture . . . . .	4
3.4	Training Process . . . . .	4
<b>4</b>	<b>Experimental results</b>	<b>5</b>
<b>5</b>	<b>Conclusion</b>	<b>6</b>

## List of Abbreviations

**ReLU**      Rectifier Linear Unit

## 0 Project Setup

For a detailed description of how to set up the project, please have a look at [https://github.com/nilskre/bomberman\\_rl/blob/master/README.md](https://github.com/nilskre/bomberman_rl/blob/master/README.md).

## 1 Introduction

## 2 Fundamentals and Related Work

### 2.1 From Language Models to modeling Evolution Theory

### 2.2 GISAID EpiFlu

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

### 2.4 Domain-Specific Methodologies to create Evolutionary Datasets for Mutation Prediction

### 2.5 Sequence2Sequence Models based on Long Short-Term Memory

- Covid-Paper: <https://www.hindawi.com/journals/mpe/2021/9980347/>
- LSTM: [https://www.researchgate.net/publication/13853244\\_Long\\_Short-term\\_Memory](https://www.researchgate.net/publication/13853244_Long_Short-term_Memory)
- Seq2Seq: <https://arxiv.org/abs/1409.3215>

### 2.6 Applying Generative Adversarial Networks

- Covid-Paper: <https://arxiv.org/pdf/2008.11790.pdf>

### 2.7 Transformer and Attention Mechanism

- Improvement: <https://arxiv.org/abs/1706.03762>

### 2.8 Other Techniques

- NNs/SVMs: <https://bsb-urasipjournals.springeropen.com/articles/10.1186/s13637-016-0042-0>
- BiLSTM: <https://science.sciencemag.org/content/371/6526/284>

## **3 Approach**

### **3.1 Dataset Creation**

### **3.2 Data Preprocessing**

### **3.3 Model Architecture**

### **3.4 Training Process**



## 4 Experimental results

## 5 Conclusion