



**FACULTY  
OF MATHEMATICS  
AND PHYSICS**  
Charles University

**MASTER THESIS**

Bc. Lukáš Chaloupský

**Automatic generation of medical  
reports from chest X-rays**

Institute of Formal and Applied Linguistics

Supervisor of the master thesis: Mgr. Rudolf Rosa, Ph.D.

Study programme: Computer Science

Study branch: Software and Data Engineering

Prague 2022

I declare that I carried out this master thesis independently, and only with the cited sources, literature and other professional sources. It has not been used to obtain another or the same degree.

I understand that my work relates to the rights and obligations under the Act No. 121/2000 Sb., the Copyright Act, as amended, in particular the fact that the Charles University has the right to conclude a license agreement on the use of this work as a school work pursuant to Section 60 subsection 1 of the Copyright Act.

In ..... date .....  
Author's signature

First of all, I would like to thank my supervisor Mgr. Rudolf Rosa, Ph.D. for all his time, guidance and valuable advices he gave me while working on this thesis. I would also like to thank my parents for their unlimited support and patience during my studies.

Title: Automatic generation of medical reports from chest X-rays

Author: Bc. Lukáš Chaloupský

Institute: Institute of Formal and Applied Linguistics

Supervisor: Mgr. Rudolf Rosa, Ph.D., Institute of Formal and Applied Linguistics

Abstract: This thesis deals with the problem of automatic generation of medical reports in the Czech language based on the input chest X-ray images using deep neural networks. The first part deals with the analysis of problem itself including comparison of existing solutions from several common points of view. In order to interpret medical images in the Czech language we present a fine-tuned a Czech GPT2 model specialized on medical texts based on the original pre-trained English GPT2 model along with its evaluation. In the second part the created Czech GPT2 is used for training neural network model for generating medical reports. The training was conducted on freely available data along with data pre-processing and their adjustment for the Czech language. Furthermore the model results are discussed and evaluated using standard metrics for natural language processing to determine the performance.

Keywords: natural language processing, image captioning, x-ray, medical report generation

# Contents

<b>Introduction</b>	<b>3</b>
<b>1 Problem Analysis</b>	<b>5</b>
1.1 Methods of generation . . . . .	5
1.1.1 Encoder . . . . .	6
1.1.1.1 Convolutional neural networks . . . . .	6
1.1.1.2 Transformers . . . . .	6
1.1.2 Decoder . . . . .	6
1.1.2.1 Recurrent neural networks . . . . .	7
1.1.2.2 GPT2 . . . . .	7
1.2 Data . . . . .	8
1.2.1 Existing datasets . . . . .	8
1.2.1.1 Indiana University chest X-ray . . . . .	9
1.2.1.2 MIMIC-CXR v2.0.0 . . . . .	9
1.2.2 Czech data . . . . .	11
1.2.3 Translators . . . . .	12
1.2.3.1 DeepL . . . . .	12
1.2.3.2 Google Translate . . . . .	12
1.2.3.3 CUBBITT . . . . .	12
1.3 Related work . . . . .	13
<b>2 Design</b>	<b>16</b>
2.1 Our approach . . . . .	16
2.2 Czech GPT-2 . . . . .	17
2.2.1 Data . . . . .	17
2.2.1.1 General . . . . .	18
2.2.1.2 Medical . . . . .	19
2.2.2 Training . . . . .	19
2.3 Medical dataset translation . . . . .	23
2.3.1 Translator choice . . . . .	23
2.3.2 Preprocessing . . . . .	23
<b>3 Experiments</b>	<b>26</b>
3.1 Environment . . . . .	26
3.2 Czech GPT-2 . . . . .	26
3.3 Dataset translation . . . . .	26
3.4 Medical report generation model . . . . .	26
<b>4 Evaluation</b>	<b>27</b>
4.1 Experiments . . . . .	27
4.2 Automatic evaluation . . . . .	27
4.2.1 Metrics . . . . .	27
4.2.2 Results . . . . .	27
4.3 Manual evaluation . . . . .	27
4.3.1 Method . . . . .	27

4.3.2 Results . . . . .	27
4.4 Examples . . . . .	27
<b>Conclusion</b>	<b>28</b>
<b>Bibliography</b>	<b>29</b>
<b>List of Figures</b>	<b>33</b>
<b>List of Tables</b>	<b>34</b>
<b>List of Abbreviations</b>	<b>35</b>
<b>A Attachments</b>	<b>36</b>
A.1 First Attachment . . . . .	36

# Introduction

In hospital, inspecting the X-rays and writing a corresponding medical reports is a hard work that requires experienced specialized doctors, of which there are not many. A great deal of people visit hospitals daily and X-rays are taken for many of them. Automatic interpretation of X-ray image has a great potential to improve health care and it could be particularly helpful to doctors in order to distinguish serious cases from the ordinary ones and overall accelerate and improve their work.

Automatic generation of radiology reports is a subset of a general problem called Image Captioning, i.e. generation of overall textual captions to input images. Image captioning is a combination of Natural Language Processing and Computer Vision areas, experiencing a lot of progress in the last years. Most often the Image Captioning problem is solved using Deep Learning techniques. The specificity of this subset is that we do not want to generate just a general caption of the image, but the exact description of all findings contained in the given medical image. There were done multiple studies for this task in other languages but none in the Czech language.

Deep learning by its very nature has wide range of uses in a medical sector as it can capture complex relations in any kind of data with excellent performance results. Nevertheless in the medical environment the accuracy of predictions is crucial in order to determine the final diagnosis. Therefore, we should not consider the models as such as something that is unmistakably true, but as an auxiliary tool that should help doctors to examine X-rays.

Inasmuch as it is not so challenging to detect fractures on the limbs, this area is less interesting than others which have a variety of diverse possible problems. One of these areas is chest for which there exists multiple freely accessible datasets containing full textual medical reports. However, all these available datasets have one common downside, they are not in the Czech language. The natural question arises, where do we obtain these much needed data? We have to face and solve this core problem in our thesis.

## Goals

First of all, we will take a closer look at the problem itself. This includes breaking down the problem and analyzing all its parts individually together with presenting possible existing alternatives for each part.

Our first goal is to fine-tune a language model directly for the Czech language. The language model will be specialized directly to medical texts in order to capture the essence of the problem. However, ahead of the medical specialization, we want to fine-tune a general czech language model. Fine-tuning will be based on the original English GPT2 model presented in Radford et al. [2019].

Finally, we want to utilize our fine-tuned language model for training neural network model interpreting chest X-rays images and generating corresponding medical textual reports to them in the Czech language. This section also involves the overall data preparation directly for the Czech language. In addition, the training will be done in multiple setups. All possibilities will be evaluated with the purpose of determination of their final performance.

## Thesis structure

In the very first chapter we present a detailed description of our problem. Every aspect of our problem is introduced and all existing solutions or possibilities are discussed with their pros and cons. Moreover we introduce there some of the important related works.

Following chapter is dealing with the design of solution to our problem, with all reasonings and decisions made. This includes not only the final neural network model, but also the language model fine-tuning and data preparation.

All experiments done with our models take their part in the third chapter. Describing all used scripts and different setups together with data variations.

Whole fourth chapter is then dedicated to evaluation of experiments done in the preceding part. Furthermore our models will be compared to the performance of other existing solutions.

Finally, in epilog we discuss what we have accomplished in the thesis, what the resulting consequences are and what the future possibilities are.



# 1. Problem Analysis

This chapter deals with the overall analysis of the problem itself. In the very beginning we present the definition of the problem. Every aspect of the problem is further discussed in detail along with a comparison of possible solutions. Moreover, the next section of the chapter describes data we work with and their alternatives. The final part of this chapter presents some of the important related works.

## 1.1 Methods of generation

First of all, we will present possible approaches for image caption generation. Inasmuch as almost all solutions for image captioning based on neural networks are using an encoder-decoder architecture as their core with further minor or major adjustments we will discuss just this architecture. Encoder serves for the purposes of extracting visual features from the input images into intermediate vector representations for the decoder. The decoder is subsequently fed with these visual feature vectors and decodes them token by token into the natural language text.

Moreover, over the past years the encoder-decoder architecture has been enhanced using the attention mechanism. The problem of encoder-decoder architecture without any other mechanism is that all the information about the image have to be encoded globally inside a single vector. However not all the information about the image is relevant for the caption generation. For this reason the attention mechanism was introduced. Instead of a single vector, the set of feature vectors representing the spatial information of the input image is used. This allows us to dynamically focus only on specific areas during the generation. The overall nature of the mechanism is derived from the way our brains concentrate on the images.

Attention mechanism for image captioning was first introduced in the Xu et al. [2015] work using the Bahdanau attention proposed in Bahdanau et al. [2014]. After extractions of visual feature vectors  $h_j$  from image, we want to assign a weight  $\alpha_{ij}$  to each of this vectors indicating the relevance of the image position  $j$  when generating output at position  $i$ . This results in the context vector  $c_i$ , a dynamic weighed combination of  $h_j$  features, that is presented as another input to the decoder. All attention calculations are defined as follows:

$$c_i = \sum_{j=1}^t \alpha_{ij} h_j \quad (1.1)$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^t \exp(e_{ik})} \quad (1.2)$$

$$e_{ij} = a(s_{i-1}, h_j) \quad (1.3)$$

where  $s_{i-1}$  is a previous hidden state and  $a(s_{i-1}, h_j)$  is an attention model.

### 1.1.1 Encoder

Encoder is responsible for extracting high-level visual features of the input image into one or more visual feature vectors. Moreover, the encoder can include also additional parts for another independent purposes, e.g. detection of objects in the image and their classification, which may be further passed to the decoder as a separate or an extra input. We will describe some of the most used neural network architectures utilized as an encoder.

#### 1.1.1.1 Convolutional neural networks

Convolutional neural networks (CNN)<sup>1</sup> are the most leveraged type of neural networks for the purposes of image analysis. They are applying 2D convolution filters for all positions in the image with the shared weights allowing to detect similar patterns independently on the position. Each layer of filters will down-size the dimensions of the previous layer and increases the number of filters. By gradually applying filters, the network is learning more complex patterns.

Instead of training convolutional neural networks from scratch, already pre-trained CNN models are often used and possibly further fine-tuned to a specific area. These models have been trained on a vast amount of data thanks to which they learned to recognize a relevant features in the lower layers applicable to any kind of images. In contrast, training from scratch is costly and takes a much longer time. Among the used pre-trained CNN model architectures are for example ResNet[He et al., 2016], VGGnet[Simonyan and Zisserman, 2014], DenseNet[Huang et al., 2017] or EfficientNet[Tan and Le, 2019].

#### 1.1.1.2 Transformers

With the Vaswani et al. [2017], a new neural network architecture called transformers was introduced. The transformers are using multi-head self-attention mechanism to compute the relations between the elements in the sequences. Although the transformers were originally designated for NLP tasks, they can be also utilized for other domains like images due to their robustness. Vision transformers (ViT), presented in Dosovitskiy et al. [2020], is an adaptation of transformers encoder for image classification without any CNNs. The core of the encoder is identical to the one in the original transformer. Nevertheless, the input image is represented as a sequence of patches for which their embeddings are computed together with the positional encodings as the input, as we can see in the Figure 1.1.

### 1.1.2 Decoder

The second main part of the network is the decoder which serves as a language model for generating corresponding captions to input images. During the generation process the visual features of the image from the encoder and the already generated text are taken into account for the prediction of the next token or word. In the process of generation, the attention mechanism (described above in

---

<sup>1</sup>[https://en.wikipedia.org/wiki/Convolutional\\_neural\\_network](https://en.wikipedia.org/wiki/Convolutional_neural_network)

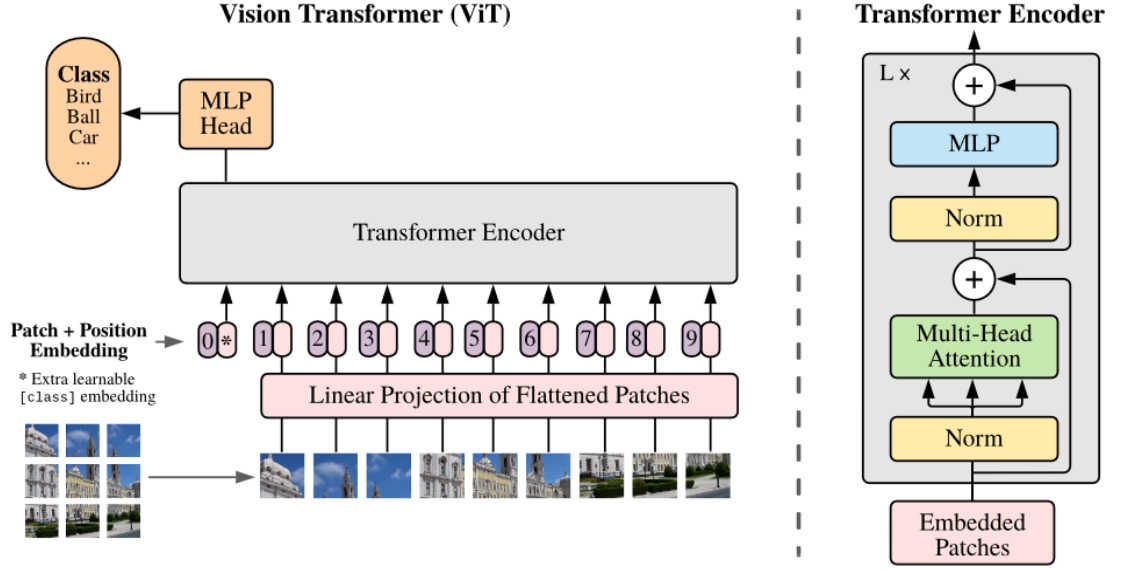


Figure 1.1: Visual Transformer architecture from Dosovitskiy et al. [2020].

the Chapter 1.1) is almost always incorporated in order to focus only on the substantial particular areas. The following subsections presents common approaches used as the backbones for the decoder part of the network.

#### 1.1.2.1 Recurrent neural networks

Recurrent neural networks (RNN)<sup>2</sup> are a category of neural networks designated for sequence data processing. Context of the previous part of the sequence can influence the subsequent elements due to the RNN cell's internal memory (hidden state). Each time step RNN cell computes its activation function from current input and hidden state producing updated hidden state as output. For language modeling tasks, the last generated token or word is given to the RNN as the next input until the entire text is produced. The autoregressiveness of the RNNs is the reason why they are ideal for use as a language model. The two most used types of RNN cells are LSTM[Hochreiter and Schmidhuber, 1997] and GRU[Cho et al., 2014], however the LSTM cells are more utilized as the decoder because they can remember longer sequences. Moreover, we can combine them in a hierarchical manner to capture more complex structures in the generated text. Nevertheless, the downside of the RNNs is the training time due to their sequential nature. The whole sentence must pass through RNN token by token and cannot be parallelized.

#### 1.1.2.2 GPT2

Just as in the case of encoder, with the advent of transformers presented in the Vaswani et al. [2017] paper, the deoder part of transformers started to be used as the language model for image captioning task for their great results in the natural language processing (NLP)<sup>3</sup> tasks. One of the advantages of the

<sup>2</sup>[https://en.wikipedia.org/wiki/Recurrent\\_neural\\_network](https://en.wikipedia.org/wiki/Recurrent_neural_network)

<sup>3</sup>[https://en.wikipedia.org/wiki/Natural\\_language\\_processing](https://en.wikipedia.org/wiki/Natural_language_processing)

transformers against the previously used recurrent neural networks is the loss of the need for sequential processing during the training. Due to the fact that the processing of the whole caption can be done in parallel, the entire training process is significantly accelerated. Moreover, the transformers can capture longer ranges dependencies in the text, as they process the sentence as a whole instead of sequentially by words.

One of the state-of-the-art autoregressive language models using transformers as their backbone is OpenAI GPT-2 model from the Radford et al. [2019] paper, which outperforms other language models on many NLP tasks. It was trained on a massive English dataset called WebText (introduced in the same article) containing a total of 40 GB of raw text. The resulting model is able to generate large coherent texts. Furthermore, it can be fine-tuned to a different domain or to a completely different language.

## 1.2 Data

In previous part we talked about possible methods of generation. Another crucial aspect we need to discuss are data, which are a basic building block of our thesis. This part focuses on the analysis of the data we used in our thesis, but also on their alternatives.

In order to solve our task and train neural network we need to get dataset containing the X-rays images along with their textual descriptions and optionally some other attributes of the examined X-rays. Moreover, the fundamental feature we need is that the data must be in the Czech language.

### 1.2.1 Existing datasets

Medical environment provides a plenty of diverse potential problems, which can be researched. As already mentioned, in this thesis we focus specifically on the X-ray images. Because it is not so hard to detect fractures on the limbs, this area is not as interesting as others. One area that is rich in its diversity is the chest. As a result, this area is explored the most and therefore there exists multiple datasets with full textual medical reports. In the following section we describe some of them.

Apart from the datasets described below, other datasets with similar type are being used with the aim of solving our task. Amongst them belong datasets such as ImageCLEFmed Caption[Rückert et al., 2022], PadChest[Bustos et al., 2020], BCIDR[Zhang et al., 2017] and PEIR Gross[Jing et al., 2017]. Moreover, except for datasets containing textual reports there exist a lot of other datasets worth mentioning containing different kind of information for each X-ray. These include, for example, CheXpert[Irvin et al., 2019], VinDr-CXR[Nguyen et al., 2020], ChestX-ray8[Wang et al., 2017] and its expanded version ChestX-ray14.

### 1.2.1.1 Indiana University chest X-ray

Indiana University chest X-Ray dataset has become a standard in the field of medical report generation, it was presented in the Demner-Fushman et al. [2015] paper. This dataset is an open source collection of pairs of chest X-rays and their corresponding semi-structured textual radiology reports, which is freely available on the web<sup>4</sup> without any additional requirements. We have a choice if we want to download just reports or images and in either PNG or DICOM format. The entire dataset consists of 7470 chest X-ray images that cover not only the frontal (PA<sup>5</sup>) view, but also the lateral (side) one. These images corresponds to a total of 3995 patient’s medical text reports.

Figure 1.1 shows an example from the Indiana University chest X-ray dataset. Each dataset pair is carefully de-identified in order to remove any personal information. The text of the report is semi-structured in up to 5 sections. The most important sections are *impression*, where the overall diagnosis is stated, *findings* section describing the details of examination and *tags* which are of two types - manual and automatic. Manual tags were annotated manually using MeSH<sup>6</sup> and RadLex<sup>7</sup> codes, automatic were encoded from the reports using the MTI indexer<sup>8</sup>. The rest of the sections are *indication* and *comparison*.

The disadvantage of this dataset is that it is relatively small. On the other hand, it is a quite clean and manually checked dataset containing also additional information about images in a form of tags described above.

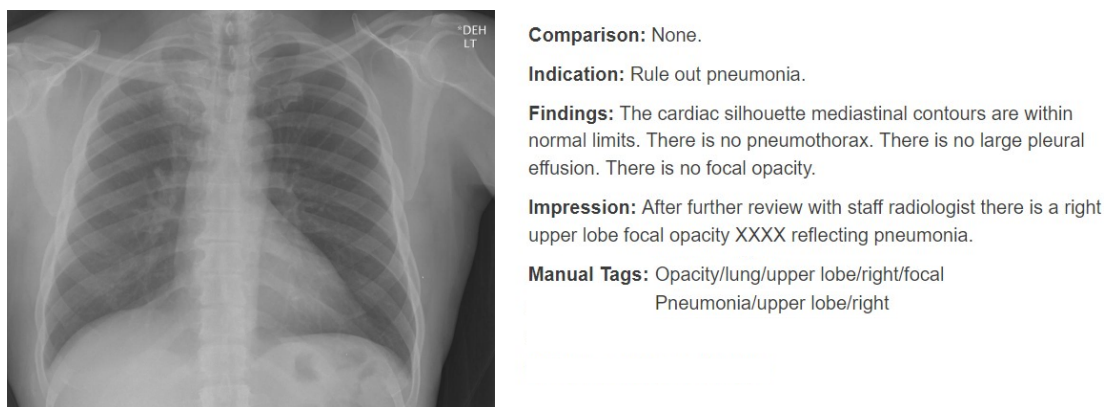


Figure 1.2: Sample from the Indiana University Chest X-ray dataset.

### 1.2.1.2 MIMIC-CXR v2.0.0

MIMIC-CXR v2.0.0 is another dataset consisting of full semi-structured medical textual reports against corresponding chest X-rays that was presented in the Johnson et al. [2019a] paper. As the previous dataset, it is openly available on

<sup>4</sup><https://openi.nlm.nih.gov/faq#collection>

<sup>5</sup>Posterior-Anterior

<sup>6</sup><https://www.nlm.nih.gov/mesh/meshhome.html>

<sup>7</sup><http://radlex.org/>

<sup>8</sup><https://lhncbc.nlm.nih.gov/ii/tools/MTI.html>

the web<sup>9</sup>. In order to get access to the dataset, we have to go through registration and verification steps. The verification phase includes completion of CITI<sup>10</sup> *Data or Specimens Only Research* course for *Human Subject Research*. Moreover we need somebody trustworthy as a reference to confirm the authenticity of our identity. After the verification we get access to all datasets in the same repository.

The dataset consists of 377,110 X-ray images in the DICOM<sup>11</sup> format connected to a total of 227,835 radiology reports for 65,379 patients. Each report is structured into multiple different sections. In order to satisfy legal requirements, entire dataset is automatically de-identified to remove any protected health information<sup>12</sup>. Similarly to the previous dataset the essential two sections of each report are *impression* and *findings*. There also exists older MIMIC-CXR-JPG<sup>13</sup> dataset, presented in the Johnson et al. [2019b] paper. This is an older version of MIMIC-CXR v2.0.0 dataset consisting of the exactly same images, only in JPG format, but each image is assigned 14 labels indicating the presence of the category in the report instead of its textual form. Each category has assigned either a 1, 0 or -1 label with the meaning *positively mentioned*, *negatively mentioned* or *uncertain*. The labels were determined from the reports utilizing the CheXpert[Irvin et al., 2019] and the NegBio[Peng et al., 2018] open-source labelers.

The advantage of this dataset is its vast number of samples. Moreover, as described above, we can get additional information in a form of categories to every image. Nevertheless the textual reports carry some noise in them in the form of grammatical mistakes and incorrect formatting. We face these issues in the Chapter 2.3.2.

---

<sup>9</sup><https://physionet.org/content/mimic-cxr/2.0.0/>

<sup>10</sup><https://about.citiprogram.org/series/human-subjects-research-hsr/>

<sup>11</sup><https://www.dicomstandard.org/>

<sup>12</sup>[https://en.wikipedia.org/wiki/Protected\\_health\\_information](https://en.wikipedia.org/wiki/Protected_health_information)

<sup>13</sup><https://physionet.org/content/mimic-cxr-jpg/2.0.0/>



#### FINAL REPORT

**EXAMINATION:** CHEST (PA AND LAT)

**INDICATION:** History: \_\_\_F with dyspnea

**TECHNIQUE:** Chest PA and lateral

**COMPARISON:** \_\_\_

**FINDINGS:** Heart size remains mild to moderately enlarged. The aorta is tortuous and diffusely calcified. Mediastinal and hilar contours are otherwise unchanged. Previous pattern of mild pulmonary edema has essentially resolved. Mild atelectasis is seen in the lung bases without focal consolidation. Blunting of the costophrenic angles bilaterally suggests trace bilateral pleural effusions, not substantially changed in the interval. No pneumothorax is present.

**IMPRESSION:** Interval resolution of previously seen mild pulmonary edema with trace bilateral pleural effusions.

Figure 1.3: Sample from the MIMIC-CXR dataset.

### 1.2.2 Czech data

All freely available datasets presented in the previous part have one common downside, namely they are not in the Czech language. As a part of elaboration of this thesis an intensive communication with real czech hospitals and other possible sources of real data took place. The goal of this communication was to create the very first open czech dataset of this kind. Processing of this kind of data would mean not only preparing the data into suitable format but also it would include proper anonymization of any personal information about the patients within the data.

However, inasmuch as the authentic patients data from hospitals are subject to strict privacy rules and we are not employees of any hospital, the institutions

decided that they cannot provide the data in any way without the conscious permission of patients given before the examination. With this result we need to find a different way how to obtain this much needed czech data.

### 1.2.3 Translators

In the previous sections we discovered that there is no dataset in the Czech language for our problem and there is no easy way how to get acces to the real data in order to build one. The only thing left is to create a new artificial dataset using an automatic translation. We will compare different freely accesible translators and choose the right one for our needs.

#### 1.2.3.1 DeepL

At the moment, DeepL<sup>14</sup> translator provides the finest available translations beating even the ones from Google Translate. Moreover, it has freely usable web application and REST API. However, the main drawback of the DeepL translator is that its REST API is highly limited - only 500 000 characters per month can be translated for free. Furthermore, any translation above this limit is costly and thus this path is not appropriate for translating large textual datasets. One way to get around this problem is to use their internal REST API used specifically for the web application, which is free to use. We investigated and implemented this potential way in our thesis and further experimented how much it can be used, but unfortunately even this internal REST API is strictly limited for only tens of consecutive<sup>15</sup> translations making it unusable for out needs.

#### 1.2.3.2 Google Translate

Google Translate<sup>16</sup> has become already de facto standard in the world of machine translation and it is the most used freely accessible language translation service in the world. In terms of quality, the translations are still great although little bit worse than those from DeepL. The web application is free of any charge and anybody can use it as much as he needs. Nevertheless, just as in the case of DeepL, their REST API services are limited and translation of anything above that limit is expensively charged. For these reasons, as in the previous case, we must find another way.

#### 1.2.3.3 CUBBITT

Machine Translation<sup>17</sup> is an extensive area of research, as a result of which there exist many other projects and academic papers nowadays. One of them is CUBBITT<sup>18</sup> translator, which was developed at our faculty. The whole system is presented and described in detail in the Popel et al. [2020] paper.

---

<sup>14</sup><https://www.deepl.com/translator>

<sup>15</sup>REST API calls are delayed from each other for some time, otherwise the service is blocked immediately

<sup>16</sup><https://translate.google.com/>

<sup>17</sup>[https://en.wikipedia.org/wiki/Machine\\_translation](https://en.wikipedia.org/wiki/Machine_translation)

<sup>18</sup><https://lindat.mff.cuni.cz/services/translation/>



CUBBITT translator provides translations which are comparable to the ones from DeepL and Google Translate services. As other mentioned translators it provides an openly available web application for machine translation. Moreover and most importantly it provides REST API that is completely unlimited in text volume and free to use without any additional charges. These are the reasons why we will utilize CUBBITT in our thesis as a translator to create our artificial dataset.

On the other hand, CUBBITT has not support for auto-correcting input text compared to above mentioned services. Moreover, there are some patterns in the text which CUBBITT cannot translate at all or translates them incorrectly. These problems complicates our situation as the data from hospitals carry some natural noise in them. We face these complications in Chapter 2.3.2.

## 1.3 Related work

The last section of this chapter is dedicated to the description and comparison to some of the related works that solves identical or similar problem as we do.

The most significant related work is Alfarghaly et al. [2021] inasmuch as we base our thesis on it. This paper focuses on the identical problem as we do, only in the English language. The proposed solution uses a pre-trained and further fine-tuned Chexnet model, presented in Rajpurkar et al. [2017], as a visual features encoder and distilGPT2<sup>19</sup> language model[Sanh et al., 2019] as a decoder which is additionally conditioned on the visual features and predicted tag’s word2vec[Mikolov et al., 2013] embeddings. For training the neural network the Indiana University chest X-ray dataset, described in more detail in Chapter 1.2.1.1, is used. Figure 1.3 shows examples of the model outputs.

Another paper making use of transformers is Chen et al. [2020]. The visual features of images are extracted using pre-trained convolutional neural network and they are further passed to the transformer encoder outputting hidden states, that are further presented to the transformer decoder for the report generation. However, the decoder architecture contains special memory module and also enhances the layer normalization. The memory module serves for memorization of text patterns which occur in the similar images inasmuch as they can further help for generating the report. We can see its effect in the Figure 1.4. Indiana University chest X-ray and MIMIC-CXR v2.0.0 datasets (see Chapter 1.2.1 for more details) were used as training data.

Yuan et al. [2019] paper proposes a hierarchical encoder-decoder architecture as we can see in the Figure 1.4 for the purpose of generating textual reports. Pairs of frontal and lateral X-ray images are used as an input to the network instead of single images, as the authors claim that the images should be complementary to each other instead of being processed independently. The RestNet-152 model pre-trained on the CheXpert[Irvin et al., 2019] dataset is utilized as the encoder with

---

<sup>19</sup><https://huggingface.co/distilgpt2>



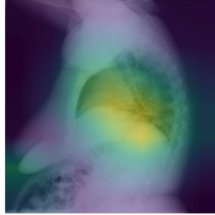
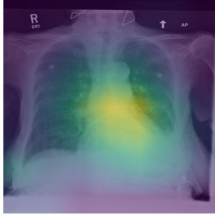
	Normal	Accurate	Missing Details	False
				
Ground Truth	Negative. The cardiomeastinal silhouette is normal in size and contour. No focal consolidation, pneumothorax or large pleural effusion. normal.	Stable cardiomegaly with mild pulmonary interstitial edema. Unchanged cardiomegaly negative for pneumothorax or focal consolidation. No large effusion. Mildly prominent interstitial opacities.	Minimal left basilar atelectasis versus infiltrate. Low lung volumes. Normal cardiomeastinal contours. Low lung volumes with minimal left basilar opacities. No pneumothorax or pleural effusions.	Prominent hiatal hernia. left basilar opacity compatible pleural effusion and atelectasis. Right pleural effusion. no pulmonary edema / overt chf identified. stable senescent mediastinal contour.
Prediction	No evidence of active Disease. The heart size and pulmonary vascularity appear within normal limits. The lungs are free of focal airspace disease. No pleural effusion or pneumothorax is seen.	Cardiomegaly and mild interstitial pulmonary edema. Moderate cardiomegaly bibasilar and perihilar interstitial opacities. No pneumothorax	No acute cardiopulmonary abnormality The lungs are clear, and without focal air space opacity The cardiomeastinal silhouette is normal in size and contour. There is no mediastinal contour, large pleural effusion, pneumothorax, or focal airspace consolidation.	Interval left subclavian central venous catheter with tip approximating the high svc. No evidence of pneumothorax. Generalized heart size and mediastinal contours appear within normal limits. Atherosclerotic changes of the aorta. Moderate degenerative changes of the thoracic spine.

Figure 1.4: Examples of generated medical reports from Alfarghaly et al. [2021].

three outputs used later - the global and local features of the images, predicted observations and medical concepts. The decoder is hierarchical LSTM decoder comprising of the two parts: *sentence decoder* and *word decoder*. The sentence decoder takes visual features and generates hidden state for each sentence, which are along with the predicted medical concepts presented to the word decoder in order to generate the report. There are many other papers using a hierarchical architecture, e.g. Huang et al. [2019].

The last related project we mention in this work is CareBot<sup>19</sup>, which is being developed concurrently with our thesis. CareBot is a czech startup founded in 2021 as a reaction to the then ongoing Covid-19 pandemic. As in our work, CareBot focuses its attention mainly on the processing of chest X-rays using neural networks. However, unlike us it does not generate the textual reports for the doctors. Instead of textual reports, it focuses on finding and classifying a total of 15 different types of individual diseases on X-ray images and their subsequent spatial localization. Moreover, they already have a support from the medical environment.

<sup>19</sup><https://www.carebot.com/>

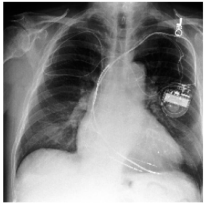
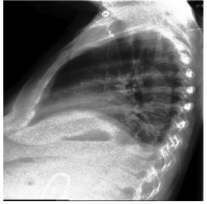
	<b>Ground-truth</b> In comparison with study of there is again enlargement of the cardiac silhouette with a pacer device in place. No definite vascular congestion raising the possibility of underlying cardiomyopathy or pleural effusion. No acute focal pneumonia. The right picc line has been removed.	<b>BASE</b> There is a left pectoral pacemaker with leads terminating in the right atrium and right ventricle. The heart is enlarged. There is no pneumothorax or pleural effusion. The lungs are clear.	<b>BASE+RM+MCLN</b> In comparison with the study of there is little change in the appearance of the pacer leads which extend to the right atrium and apex of the right ventricle. Continued enlargement of the cardiac silhouette without vascular congestion or pleural effusion. No evidence of pneumothorax.
	<b>Ground-truth</b> There are low lung volumes. Bibasilar atelectasis have minimally improved. Mild vascular congestion has minimally improved. There are no new lung abnormalities or pneumothorax. Bilateral pleural effusions are small. Right picc tip is at the cavoatrial junction.	<b>BASE</b> In comparison with the study of there is little overall change. Again there is some indistinctness of pulmonary vessels consistent with elevated pulmonary venous pressure. No evidence of acute focal pneumothorax.	<b>BASE+RM+MCLN</b> The lung volumes are low. There is a small left pleural effusion with associated atelectasis. The right lung is clear. There is no pneumothorax. The heart size is top normal. The hilar and mediastinal contours are normal. A right subclavian catheter terminate in the mid svc.

Figure 1.5: Examples of generated medical reports from Chen et al. [2020].

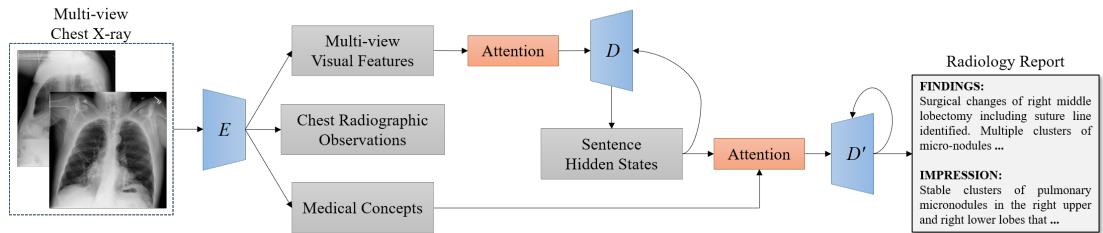


Figure 1.6: Hierarchical architecture from Yuan et al. [2019].  
 $E$  is encoder,  $D$  is sentence decoder and  $D'$  is word decoder.

## 2. Design

This chapter summarizes the overall design of our approach to generating textual medical reports for X-ray images. The problem consists of multiple independent parts we need to deal with. For each of them, we will present the fundamentals of our solution, along with description of related problematics and decisions made.

### 2.1 Our approach

As we already mentioned in the Chapter 1.3, the overall solution for the final medical report generation model is based on the Alfarghaly et al. [2021]. We have chosen this approach for multiple reasons. The main reasons to use this work as the backbone for our thesis are following:

1. In the work the state-of-the-art GPT-2 model is utilized as the language model. This gave us a great opportunity as there was none Czech GPT-2 model available at the time this thesis began.
2. The encoder is already fine-tuned to extract visual features for specific dataset.
3. All solution source code is freely accessible on the github<sup>1</sup>.

As in most works for image captioning, the architecture is encoder-decoder based with an attention mechanism. The high-level solution architecture is depicted in Figure 2.1.

The encoder part utilizes fine-tuned Chexnet[Rajpurkar et al., 2017] as its visual backbone. The Chexnet is a CNN Densenet121 model trained specifically for medical environment on the ChestX-ray14[Wang et al., 2017] dataset predicting 14 distinct disease classes from the input chest X-ray images. Input images are resized to the  $224 \times 224$  resolution as in the case of the original Chexnet model. Moreover, the encoder is further fine-tuned to predict 105 of the most common manual tags from the Indiana University chest X-ray dataset(see Chapter 1.2.1.1) for the purpose of wider range of possible semantic features. The model thus produces two types of outputs - visual features from the base CNN model and class scores from the multi-label tag prediction.

In order to obtain semantic features from the predicted classes, each predicted class score is multiplied by its corresponding word2vec[Mikolov et al., 2013] embedding, trained specifically on the biomedical texts, from the McDonald et al. [2018] and resulting in the weighted embedding matrix. Both the embedding matrix and visual features are further passed to the language model as the context vectors for the self-attention mechanism.

---

<sup>1</sup><https://github.com/omar-mohamed/GPT2-Chest-X-Ray-Report-Generation>

Instead of the original English GPT-2 the original solution used its smaller version distilGPT-2<sup>2</sup>. The distilled version differs in the number of layers as it contains only a half of the layers compared to the original model and predicts only 512 tokens instead of 1024. Nevertheless, distilled versions of GPT-2 do not usually work well for languages other than English due to their reduced capacity. For these reasons, we use the original architecture of the GPT-2 for further fine-tuning. But, our fine-tuned model will predict 512 tokens as well as the distilGPT-2. During the medical report generation training, the maximal sequence length is set to 200.

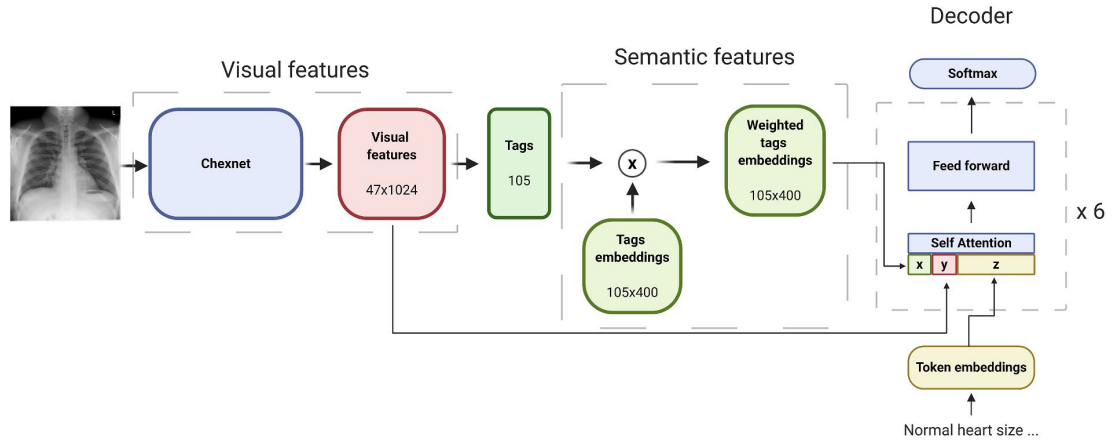


Figure 2.1: Overall architecture used in our solution proposed in Alfarghaly et al. [2021].

## 2.2 Czech GPT-2

The aim of this work is to generate medical reports in the Czech language. In the previous parts, we decided to use the GPT-2 as the language model. However at the time of the beginning of this work, no Czech GPT-2 model was freely available and thus it was essentially necessary to create one. This section describes all the steps needed for fine-tuning the small English GPT-2 to the Czech language. Respectively, we train two versions of the Czech GPT-2 model. One trained on the general Czech textual data and one specialized specifically on Czech medical texts.

### 2.2.1 Data

In this section we will describe the possible data applicable for training of both the general and medical Czech GPT-2 model together with the decision made about the final data selection and data cleaning.

<sup>2</sup><https://huggingface.co/distilgpt2>

### 2.2.1.1 General

For the training of general Czech GPT-2 we have a plenty of data options we can choose from. However, there are important properties of the data we need to satisfy. As we are transfer learning from English to Czech language we need to have sufficiently large data, so we ensure the GPT-2 will learn properly syntactical and semantical information. Moreover, the data have to be also heterogeneous enough, so the model can capture different types of information and not just for example newspaper articles from a specific area. In order to create a good enough general model we need to meet these criteria.

Several different datasets were investigated and tested for the training of the general Czech GPT-2 model.

#### Czech Wikipedia

The first data we used for training the model is the Czech Wikipedia dump<sup>2</sup>. After extraction, the total size of the dataset is approximately 800 MB of raw text. The advantages of this dataset are its easy accessibility and fairly clean data quality. On the other hand, the data are very homogeneous despite the various topics. Each article is written in the general descriptive style. Moreover the data themselves are not large enough, the trained model made many both syntactical and semantical mistakes during the text generation.

#### Balanced Czech National Corpus

Another possibility was to use a balanced version of the Czech National Corpus[Křen et al., 2021] as the original is composed mainly of journalistic articles. The balanced version tries to equalize the amount of data from each category. These categories include *journalism*, *poetry*, *prose*, *educational literature* etc. The major advantage of this dataset is its purity, the texts are syntactically correct without any undesirable non-Czech elements and written in the standard Czech language. The dataset does not have any significant downsides and the trained model understood Czech language without any significant ailments. In total, the dataset is comprised of 3,3 GB of raw text.

#### OSCAR

OSCAR, from the Ortiz Suárez et al. [2020], is a huge deduplicated multilingual corpus created from the Common Crawl corpus<sup>3</sup> providing data for 166 different languages and available directly in the huggingface datasets library<sup>4</sup>. It consists of the text scraped from websites of very different kinds and thus the data are heterogeneous enough. Moreover, its huge size, as the Czech part of the dataset occupies a total of 24 GB of raw text, is another major benefit. On the other hand, because the data are automatically scraped, they carry a noise in them. Besides that, not negligible part of the text are in the non-standard Czech

---

<sup>2</sup><https://dumps.wikimedia.org/cswiki/latest/cswiki-latest-pages-articles.xml.bz2>

<sup>3</sup><https://commoncrawl.org/>

<sup>4</sup><https://huggingface.co/datasets/oscar>

language as the data come from diverse web sources such as forums etc. Nevertheless, the disadvantages are outweighed by the huge size of the corpus and together with following filtering of the text:

1. We take only text that are at least 1200 characters long as these texts tend to be longer articles written in the standard Czech language instead of advertisements, incomplete texts etc.
2. Any text containing control character are filtered out, because the text contains generally undesirable content.

## Conclusion

We analyzed various datasets along with their overall properties. Furthermore, the advantages and disadvantages of each were discussed. As a result, we have chosen the **OSCAR** dataset due to its size and heterogeneity. The **Wikipedia** is too small and homogeneous. On the other hand the **Balanced Czech National Corpus** is heterogeneous enough, however it is an almost order of magnitude smaller than **OSCAR**.

### 2.2.1.2 Medical

Since we have a trained general Czech GPT-2 model from the previous section that already understands a Czech language, the final fine-tuning for medical environment does not require that much data. We need to specialize the model to understand the medical environment inherently. For this purpose, we use a subset of the UFAL Medical Corpus v. 1.0<sup>5</sup>. These data are further filtered to remove any inappropriate characters, lines and redundant structures. As a result, the data contain a total of 100 MB of raw medical texts. The texts comprise of general medical descriptions, articles and package leaflets for medicines.

## 2.2.2 Training

In this part, we will describe the fine-tuning process of our Czech GPT-2 models. Both the general and the medical GPT-2 models will be trained using the same process. The solution is based on the Guillou [2020] article using specific `fastai`<sup>6</sup> library providing with powerful tools using best practices aiming for training and fine-tuning neural network models. The reason why we used this article is that it has shown great results for transfer learning the English GPT-2 to another language in a short time against traditional fine-tuning. The whole process comprises of several important parts, that we will describe in the subsequent part of the text. Detailed information about the insides and experiments done are described in the following Chapters TODO.

### Training tokenizer

First thing that needs to be done in the whole training process is to train a tokenizer specifically for the Czech language. As in the case of original GPT-2

---

<sup>5</sup>[https://ufal.mff.cuni.cz/ufal\\_medical\\_corpus](https://ufal.mff.cuni.cz/ufal_medical_corpus)

<sup>6</sup><https://www.fast.ai/>

model, we use the same byte-level byte-pair-encoding<sup>7</sup> tokenizer dividing input text into tokens (a word or its part). Original size of the vocabulary is kept and set to 50257, as well as the original special token `<|endoftext|>` as indication of the beginning/end of the sequence token and the pad token. The entire prepared datasets from Chapter 2.2.1 are used for the tokenizers training.

## Initialization

As we have indicated several times, for the GPT-2 model weights initialization we will use the weights from the original English GPT-2. However, we will do a modification of the initial word token embeddings.

For the word token embedding, we will compare what tokens have the English model tokenizer and our tokenizer in common. The embeddings of these tokens will remain unchanged, as they have been already trained on a large text corpus and they already carry information, even though they come from a completely different language. The rest of the tokens will be initialized using the mean value of the English word token embeddings.

## Data preparation

Another step in the process is to prepare the dataset into suitable structure for the training. Entire dataset is loaded and pre-tokenized in advance in order to reduce the data transfer time between CPU and GPU needed. After the pre-tokenization process, all texts are further passed to a specific language modeling data object, that is responsible for preparation and handling of training and validation data. All texts are concatenated, split by the defined sequence length and formed into the batches. We use batch size of 16 as it is the maximum size we were able to fit into the GPU and sequence length of 512, because we need to have model corresponding to the one used in Alfarghaly et al. [2021] - distilGPT2 - as our main goal is to generate medical reports.

The batch size is one of the important hyper-parameters. In the following section we will discuss, how learning rates are determined and also that we will use higher learning rates. Nevertheless, in the Smith [2018] they recommend to use batch size as large as possible, because higher learning rates are regularization themselves and therefore other regularizations can be reduced.

## Finding optimal learning rate

Fine-tuning is a very fragile process in terms of learning rates. Right choice of the learning rate is essential, if we choose a very small learning rate, the model will train slowly and will tend to overfitting. On the other hand, if we choose too high learning rate, the whole process can diverge and all the progress will be lost.

Smith [2017] came with a solution to this problem. Before the training itself, we do a pre-training run in which we are trying a wide range of learning rates for which we monitor their behaviour. Starting from a very low learning rates up to

---

<sup>7</sup><https://huggingface.co/course/chapter6/5?fw=pt>



the very high learning rates, for every mini-batch we try a current learning rate, collect the resulting loss and move to the next iteration with a little higher learning rate. In the end, we plot collected losses against the corresponding learning rates.

An illustration of the result of this process can be seen in the Figure 2.2. The graph gives us an overview for which learning rates the model is still learning. General recommendation is to use a learning rate that is an order of magnitude less than the minimum as the minimum is very close to the moment of divergence. We can also see in the graph that the current implementation gives us several other significant points usable for training.

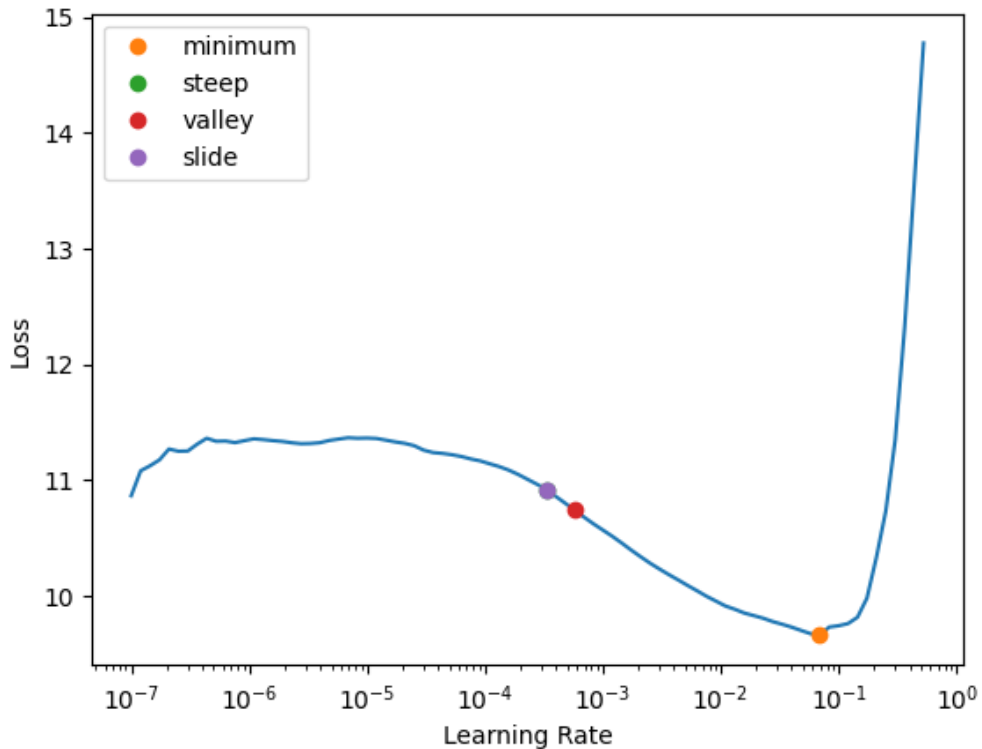


Figure 2.2: Learning finder output.  
The *Learning rate* axis is in logarithmic scale.

## Training process

We have prepared everything we need in the previous parts. However, we still have to describe the essence of our fine-tuning process. We will first describe the learning rate schedule. Fine-tuning or transfer learning is usually done using small learning rates with decay, so the general information in the already pre-trained weights are kept while the model is specialized to a selected task. The disadvantage of this method is the long training time. Instead, we use a different approach from Smith [2018] called "1cycle" policy that is parametrized by *minimal* and *maximal* learning rates. In the beginning we start with a *minimal* learning rate and linearly increase it up to the *maximum*. The second phase is then the opposite direction going down with cosine annealing, but instead of stop-

ping at the *minimum* the decrease continues down by several order of magnitudes lower as we can see in the Figure 2.3.

The intuition behind this policy is quite natural. In the beginning we start with a lower learning rates to find an optimal direction. As we are increasing the learning rate we are taking bigger steps this direction, skipping sharp local minima and preferring wide flat local minima area as shown in the Smith and Topin [2019]. The rest of the training is intended for improvement inside this area. This allows us to use higher learning rates and thus overall accelerate the entire training process. The *maximal* learning rate should be chosen according to the previous section and the *minimal* is defined as  $min = max/25$  by default.

Another major part of the training are discriminative learning rates. This method has been proposed in the Howard and Ruder [2018]. Instead of training all layers at once, we assign a different learning rates for each layer or a group of layers as each of them. This arises from the reasoning about what the different parts of the network are focusing on. Therefore, we split our model into 4 distinct groups of layers and train each of them with different learning rate.

The fine-tuning process itself is divided into two parts: training the new head for the Czech language only and training the whole model at once. In the original article, they suggest gradual unfreezing approach, where they unfreeze one more layer each time and train for one epoch. However after numerous experiments, we observed that unfreezing of the whole model and running multiple epochs at once gives better results.

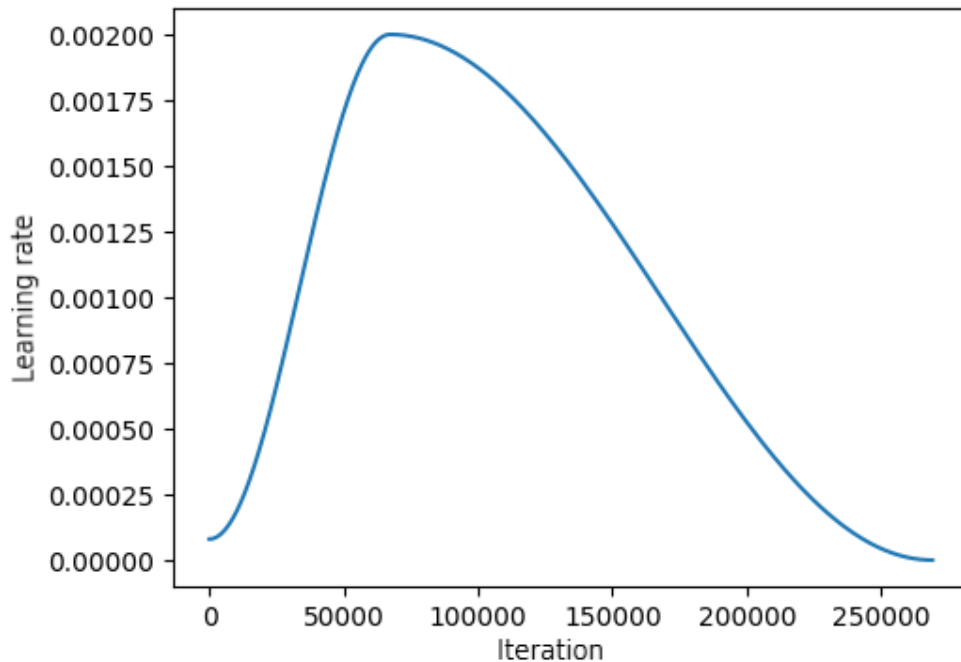


Figure 2.3: Learning rate schedule for 1cycle policy training. The graph depicts the learning rate throughout the entire training process.

## 2.3 Medical dataset translation

For machine learning in general and NLP tasks especially, the data quality is the alpha and omega of the performance of the final model. Since, as we have already mentioned in the Chapter 1.2.2, we do not have any Czech data directly for the medical examinations of X-rays, to obtain Czech data we must arrange ourselves in a different way. One potential way is to create a new artificial dataset using machine translation of existing datasets. This section discusses the required steps to build a quality dataset using translation.

### 2.3.1 Translator choice

In the previous part of the text, specifically in the Chapter 1.2.3, we already discussed all possibilities for automatic translation. Our final choice for the translator is the CUBBITT as it provides REST API unlimited in the number of requests and volume.

### 2.3.2 Preprocessing

The most important part of our machine translation process is the preprocessing of the input text. We already outlined in the Chapter 1.2.1 that the data contain some noise in them as we are dealing with reports in natural language. Moreover, we will use CUBBITT translator, that does not perform auto-correction itself and cannot translate some patterns at all as we described in the Chapter 1.2.3.3. For these purposes we incorporate preprocessing before the translation as it would be beneficial to have all texts in standardized form in order to firstly, help CUBBITT with translation to get the report correctly translated, and secondly to ensure that our model receives and process all the data in an identical report format.

Our preprocessing pipeline encompasses of following procedures. Some of them are dealing with general CUBBITT issues and other with specifics of the medical data.

#### Line starts

First of all we start with a very simple procedure. We analyze all lines of the report and delete all white characters common for all lines on both ends. The purpose of this modification is to standardize the report format, while preserving its structure.

#### Anonymous sequences

Inasmuch as the whole datasets are anonymized due to the legal reasons and privacy protection, the reports contain “anonymous sequences”, such as “XXXX” or “\_\_\_”, denoting places with original private information about patients. However, these sequences can be attached to surrounding words forming undesirable words. As we already said, the CUBBITT does not auto-correct its input automatically, so these inputs will not be handled in any special way and therefore

not translated. For this reason, we separate these sequences to form independent words.

## Units

Subsequent form of correction that we perform is the separation of numbers and units attached to them. In addition, this also includes general cases, where the number and subsequent word are glued together, while keeping specific medical terms with a similar structure. This is associated with the following step, as the units will not be true-cased properly without this procedure.

## True-casing

The most important part of the whole preprocessing pipeline is true-casing of the input text. This adjustment is necessary for two essential reasons. CUBBITT has problems with translation of any uppercase texts in general. Capturing the true-case of a text is a complex problem requiring either a large statistical language dictionary or a trained model in order to properly determine the case.

As medical reports are very specific area, the existing solutions for general text true-casing is inapplicable. Medical reports contain a lot of abbreviations and acronyms, which can be often confused with ordinary english words. Training the model for medical true-casing requires even more specific data for a certain domain, because in different contexts the common words can be treated differently. Moreover, obtaining flawless data to cover the entire specific domain is a challenging task. For these reasons, we chose the way of a statistical dictionary. Before the translation of the dataset begins, we create a statistical dictionary from the whole dataset containing the most often used form of every word. We also count with some exceptions, such as headings, that in some datasets can be in uppercase only.

Using the created dictionary, we deal with all uppercase words to assign them the proper form. The results of the true-casing preprocessing are demonstrated in the following examples:

- (1a) “EXAMINATION: CHEST (PORTABLE AP)” →  
“PŘEZKOUŠENÍ: CHEST (PORTABLE AP)”
- (1b) “Examination: Chest (portable AP)” →  
“Vyšetření: Hrudník (přenosný AP)”
- (2a) “SMALL RIGHT PLEURAL ABNORMALITY” →  
“MALÉ PRÁVO PLEURÁLNÍ ABNORMALITY”
- (2b) “Small right pleural abnormality” →  
“Malá pravá pleurální abnormalita”

## Paragraphs structure

In some of the medical reports the section headings and corresponding texts do not begin on the same lines. We adjust these situations to a form where

each heading and the text belonging to it always starts on the same line, for two reasons. Firstly, we want to normalize the report structure in general and secondly, we want to move the section content as close as possible to the heading, so the translator and even the final model have the context close to each other.

## **Capitalization**

Another part of the preprocessing pipeline is a simple capitalization of each heading and each sentence in the report. This text capitalization process helps CUBBITT not only in the case of medical data, but in general during the translation process of some texts to better understand the boundaries between sentences.

## **Time**

One of the patterns that CUBBITT does not recognize nor auto-corrects itself in general are times. This applies both, to the specification of hours and minutes, and to the part of the day specification. If any of these parts are incorrectly formatted, the time will be translated incorrectly or not translated at all, and thus we would lose some information or it could damage the fluency of the translated text. For these reasons, we apply preprocessing to normalize all times. We can see the difference in the following examples:

- “Chest radiograph at 1045PM” → “Rentgen hrudníku v 1045PM”
- “Chest radiograph at 10:45 PM” → “Rentgen hrudníku ve 22:45”

## **White spaces**

After all previous procedures, we apply one last very simple final modification, namely, we squash all the whitespace characters inside each line into a single space, while maintaining the format from the very first preprocessing step described above. This step is performed only for normalization purposes.

## **Lowercasing**

The last preprocessing procedure we will mention is lowercasing of the whole text followed by capitalizing the first word of each sentence. This is a separate procedure that was used in the earlier phases of elaboration of this work as CUBBITT has problem with uppercased word as we already mentioned.

## 3. Experiments

Úvod

### 3.1 Environment

Popsat prostředí na kterém jsme trénovali - it4i cluster a aic cluster

### 3.2 Czech GPT-2

Popis toho tréninku - odkaz na předchozí kapitolu. Výsledky z trénování \* Normal Czech GPT-2 \* Medical Czech GPT-2 Ukázky výsledků z obou Popis scriptů

### 3.3 Dataset translation

Popsat, jak vypadá ten script, jak funguje a že jsme nakonec překládali jen open-i + Výsledky

### 3.4 Medical report generation model

Popsat, že jsme použili data z Open-I jako v původním článku, jelikož nám jde hlavně o proof-of-concept + Mimic nemá stejné tagy -> muselo by se přetrénovat celý CNN backbone. + Výsledky + Popis toho, jak jsme upravili data do stejného formátu jako v článku, popis těch skriptů je na githubu toho článku

## 4. Evaluation

### 4.1 Experiments

### 4.2 Automatic evaluation

#### 4.2.1 Metrics

#### 4.2.2 Results

### 4.3 Manual evaluation

#### 4.3.1 Method

#### 4.3.2 Results

### 4.4 Examples

# Conclusion



# Bibliography

- Omar Alfarghaly, Rana Khaled, Abeer Elkorany, Maha Helal, and Aly Fahmy. Automated radiology report generation using conditioned transformers. *Informatics in Medicine Unlocked*, 24:100557, 2021.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- Aurelia Bustos, Antonio Pertusa, Jose-Maria Salinas, and Maria de la Iglesia-Vayá. Padchest: A large chest x-ray image dataset with multi-label annotated reports. *Medical image analysis*, 66:101797, 2020.
- Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan. Generating radiology reports via memory-driven transformer. *arXiv preprint arXiv:2010.16056*, 2020.
- Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- Dina Demner-Fushman, Marc D. Kohli, Marc B. Rosenman, Sonya E. Shooshan, Laritza Rodriguez, Sameer Antani, George R. Thoma, and Clement J. McDonald. Preparing a collection of radiology examinations for distribution and retrieval. *Journal of the American Medical Informatics Association*, 23(2):304–310, 07 2015. ISSN 1067-5027. doi: 10.1093/jamia/ocv080. URL <https://doi.org/10.1093/jamia/ocv080>.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- P Guillou. Faster than training from scratch—fine-tuning the english gpt-2 in any language with hugging face and fastai v2 (practical case with portuguese), 2020.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. *arXiv preprint arXiv:1801.06146*, 2018.
- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.

- Xin Huang, Fengqi Yan, Wei Xu, and Maozhen Li. Multi-attention and incorporating background information model for chest x-ray image report generation. *IEEE Access*, 7:154808–154817, 2019.
- Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghighi, Robyn Ball, Katie Shpan-skaya, et al. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 590–597, 2019.
- Baoyu Jing, Pengtao Xie, and Eric Xing. On the automatic generation of medical imaging reports. *arXiv preprint arXiv:1711.08195*, 2017.
- Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Roger G Mark, and Steven Horng. Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific data*, 6(1):1–8, 2019a.
- Alistair EW Johnson, Tom J Pollard, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Yifan Peng, Zhiyong Lu, Roger G Mark, Seth J Berkowitz, and Steven Horng. Mimic-cxr-jpg, a large publicly available database of labeled chest radiographs. *arXiv preprint arXiv:1901.07042*, 2019b.
- Michal Křen, Václav Cvrček, Jan Henyš, Milena Hnátková, Tomáš Jelínek, Jan Koček, Dominika Kovářiková, Jan Křivan, Jiří Milička, Vladimír Petkevič, Pavel Procházka, Hana Skoumalová, Jana Šindlerová, and Michal Škrabal. SYN v9: large corpus of written czech, 2021. URL <http://hdl.handle.net/11234/1-4635>. LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.
- Ryan McDonald, Georgios-Ioannis Brokos, and Ion Androutsopoulos. Deep relevance ranking using enhanced document-query interactions. *arXiv preprint arXiv:1809.01682*, 2018.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26, 2013.
- Ha Q Nguyen, Khanh Lam, Linh T Le, Hieu H Pham, Dat Q Tran, Dung B Nguyen, Dung D Le, Chi M Pham, Hang TT Tong, Diep H Dinh, et al. Vindr-cxr: An open dataset of chest x-rays with radiologist’s annotations. *arXiv preprint arXiv:2012.15029*, 2020.
- Pedro Javier Ortiz Su’arez, Laurent Romary, and Benoit Sagot. A monolingual approach to contextualized word embeddings for mid-resource languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1703–1714, Online, July 2020. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/2020.acl-main.156>.

- Yifan Peng, Xiaosong Wang, Le Lu, Mohammadhadi Bagheri, Ronald Summers, and Zhiyong Lu. Negbio: a high-performance tool for negation and uncertainty detection in radiology reports. *AMIA Summits on Translational Science Proceedings*, 2018:188, 2018.
- Martin Popel, Marketa Tomkova, Jakub Tomek, Łukasz Kaiser, Jakob Uszkoreit, Ondřej Bojar, and Zdeněk Žabokrtský. Transforming machine translation: a deep learning system reaches news translation quality comparable to human professionals. *Nature Communications*, 11(4381):1–15, 2020. ISSN 2041-1723. doi: 10.1038/s41467-020-18073-9. URL <https://www.nature.com/articles/s41467-020-18073-9>.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.
- Pranav Rajpurkar, Jeremy Irvin, Kaylie Zhu, Brandon Yang, Hershel Mehta, Tony Duan, Daisy Ding, Aarti Bagul, Curtis Langlotz, Katie Shpanskaya, et al. Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning. *arXiv preprint arXiv:1711.05225*, 2017.
- Johannes Rückert, Asma Ben Abacha, Alba García Seco de Herrera, Louise Bloch, Raphael Brüngel, Ahmad Idrissi-Yaghir, Henning Schäfer, Henning Müller, and Christoph M. Friedrich. Overview of ImageCLEFmedical 2022 – caption prediction and concept detection. In *CLEF2022 Working Notes*, CEUR Workshop Proceedings, Bologna, Italy, September 5-8 2022. CEUR-WS.org.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Leslie N Smith. Cyclical learning rates for training neural networks. In *2017 IEEE winter conference on applications of computer vision (WACV)*, pages 464–472. IEEE, 2017.
- Leslie N Smith. A disciplined approach to neural network hyper-parameters: Part 1–learning rate, batch size, momentum, and weight decay. *arXiv preprint arXiv:1803.09820*, 2018.
- Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of neural networks using large learning rates. In *Artificial intelligence and machine learning for multi-domain operations applications*, volume 11006, pages 369–386. SPIE, 2019.
- Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

- Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M Summers. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2097–2106, 2017.
- Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pages 2048–2057. PMLR, 2015.
- Jianbo Yuan, Haofu Liao, Rui Luo, and Jiebo Luo. Automatic radiology report generation based on multi-view image fusion and medical concept enrichment. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 721–729. Springer, 2019.
- Zizhao Zhang, Yuanpu Xie, Fuyong Xing, Mason McGough, and Lin Yang. Md-net: A semantically and visually interpretable medical image diagnosis network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6428–6436, 2017.

# List of Figures

1.1	Visual Transformer architecture from Dosovitskiy et al. [2020]. . .	7
1.2	Sample from the Indiana University Chest X-ray dataset. . . . .	9
1.3	Sample from the MIMIC-CXR dataset. . . . .	11
1.4	Examples of generated medical reports from Alfarghaly et al. [2021].	14
1.5	Examples of generated medical reports from Chen et al. [2020]. . .	15
1.6	Hierarchical architecture from Yuan et al. [2019]. . . . .	15
2.1	Overall architecture used in our solution proposed in Alfarghaly et al. [2021]. . . . .	17
2.2	Learning finder output. . . . .	21
2.3	Learning rate schedule for 1cycle policy training. . . . .	22

# List of Tables

# List of Abbreviations

GPT-2 – Generative Pre-trained Transformer 2

CNN – Convolutional Neural Network

RNN – Recurrent Neural Network

LSTM – Long Short-Term Memory

GRU – Gated Recurrent Unit

ViT – Vision Transformer

NLP – Natural Language Processing

PA – Posterior-Anterior

DICOM – Digital Imaging and Communications in Medicine

MIMIC – Medical Information Mart for Intensive Care

CUBBITT – Charles University Block-Backtranslation-Improved Transformer Translation

# A. Attachments

## A.1 First Attachment