



**FACULTY
OF MATHEMATICS
AND PHYSICS**
Charles University

MASTER THESIS

Bohdan Ihnatchenko

Multi-Target Machine Translation

Institute of Formal and Applied Linguistics

Supervisor of the master thesis: doc. RNDr. Bojar Ondřej, Ph.D.

Study programme: Computer Science

Study branch: Artificial Intelligence

Prague 2020

This is not a part of the electronic version of the thesis, do not scan!

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

Dedication.

Title: Multi-Target Machine Translation

Author: Bohdan Ihnatchenko

Institute: Institute of Formal and Applied Linguistics

Supervisor: doc. RNDr. Bojar Ondřej, Ph.D., Institute of Formal and Applied Linguistics

Abstract: Abstract.

Keywords: Machine translation words

Contents

Introduction	3
1 Background	4
1.1 XXX TODO: History of machine translation	4
1.2 XXX TODO: Transformer model	4
1.3 XXX TODO: Translation evaluation	5
1.3.1 History	5
1.3.2 XXX TODO: BLEU - bilingual evaluation understudy . .	5
1.4 Multi-target machine translation	6
1.4.1 Multi-lingual machine translation	6
1.4.2 Massively multi-lingual machine translation with complete sharing	7
1.5 Conclusion	7
2 Experiment setup	10
2.1 XXX TODO: Questions and constraints	10
2.2 Experiments	10
2.2.1 Starting point	10
2.2.2 Proposed experiments	11
2.3 Dataset(s)	11
2.3.1 TO EDIT: English to 36 languages	11
2.3.2 XXX TODO: UN parallel corpus: English to 5 languages .	13
2.4 Method	13
2.4.1 TO EDIT: Training	13
2.4.2 TO EDIT: Validation	14
2.4.3 XXX TODO: Finishing the training	14
2.4.4 TO EDIT: Testing	16
2.4.5 XXX TODO: Analysis	17
2.5 Training tools	17
2.5.1 Tool kits	17
2.5.2 Computational cluster	18
2.5.3 Inspecting the training process	19
2.5.4 XXX TODO: Model settings	19
3 Bilingual and multi-lingual baselines	23
3.1 XXX TODO: Bilingual baseline	23
3.2 XXX TODO: Multilingual baseline	24
3.3 Expected results	24
3.4 Performance drop on massively multilingual setup	25
3.5 Performance decrease on richer data sets	25
4 Group by language groups	27
4.1 Language groups	27
4.1.1 Germanic group	27
4.1.2 Slavic with cyrillic script	30

5 Discussion	32
5.1 Results	32
5.2 Further work	32
Conclusion	33
Bibliography	34
List of Figures	36
List of Tables	37
List of abbreviations	38
Glossary	39
A Attachments	40
A.1 Additional tables	41
A.1.1 Bilingual results	41

Introduction

With increasing availability of computational resources and enormous amount of publicly available corpora it is now possible to obtain a machine translation (MT) system which produces translations of acceptable quality. But in the use cases similar to conferences, where one speech is translated into multiple target languages, the same amount of models needs to be deployed. Another option is to use multilingual MT system for all needed languages together, which may lead to a decreased quality of translations.

1. Background

1.1 XXX TODO: History of machine translation

XXX is it ok to quote the whole paragraph with short history? e.g. Han [2016]

1.2 XXX TODO: Transformer model

Introduced in Vaswani et al. [2017] Transformer model is used as a base for numerous state-of-the-art systems as can be seen for example in WMT18 [Bojar et al., 2018] and WMT19 [Barrault et al., 2019] results.

Prior to invention of the *Transformer* model, recurrent neural network (RNN) and convolutional neural network (CNN) architectures were used to encode source side of the sentence pair and to decode it into the target sentence. Various window lengths in CNN architectures allowed to capture long range relations as well as short range ones; still the range was limited by the maximum window length. In RNN-like architectures long short-term memory (LSTM) and gated recurrent unit (GRU) cells were used, as their structure allowed to pass the internal state on longer distances due to selective forgetting.

Transformer model uses *self attention* mechanism to encode contextual information in each word position. *Position encoding* allows passing the position information without explicit sequential connections as in RNN. As was stated by *Transformer*'s authors, there are three main points why self-attention mechanism should be preferred:

- total computational complexity per layer;
- the amount of computation that can be parallelized;
- the path length between long-range dependencies in the network.

Layer type	Complexity per layer	Sequential operations	Maximum path length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

Table 1.1: **Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types.** n is the sequence length, d is the representation dimension, k is the kernel size of convolutions and r the size of the neighborhood in restricted self-attention.

1.3 XXX TODO: Translation evaluation

1.3.1 History

In 1966 first machine translation evaluation methods were proposed by the Automatic Language Processing Advisory Committee (ALPAC). The proposed metrics were "intelligibility" and "fidelity"[ALPAC, 1966, p 67]. Metrics were considered independent and the evaluation was meant to be conducted by two independent groups of raters. *Intelligibility* was measured without reference to the original by group of raters fluent in target language and not familiar with the source language. The raters were comparing the "informativeness" of evaluated translation with carefully prepared reference translation. The metric scale was from 1 (hopelessly unintelligible) to 9 (perfectly clear and intelligible). *Fidelity* to the sense of the original text was measured by another group of raters, which were native speakers of the target language and highly proficient in the source language. It was measured on a scale of 0-9 and showed how much additional information was added in the evaluated translation comparing to the reference translation.

Advanced Research Projects Agency (ARPA) XXX TODO: from White 1999 and Church 1993 comprehension evaluation, quality panel evaluation, and evaluation based on adequacy and fluency

Automatic evaluation

- Metric should correlate with human judgments;
- Works with texts of different domains
- Banerjee et al. (2005) correlation, sensitivity, consistency, reliability and generality.

1.3.2 XXX TODO: BLEU - bilingual evaluation understudy

In Papineni et al. [2002] was introduced a novel method of automatic machine translation evaluation - bilingual evaluation understudy (BLEU). Its advantages are high speed and low cost of evaluation, language independence and high correlation with judgements of highly skilled human raters evaluations.

Shortly, BLEU score incorporates modified n-gram precision scores corrected by brevity penalty, which ensures the produced translation length is close to the reference one. BLEU score is computed for the whole test corpus.

Modified *n*-gram precision score

The main element of the metric is the *precision* measure. To compute precision, the number of candidate translation words (unigrams) that are present in any reference translation is divided by the total number of words in the candidate translation. This approach leads to overrating candidate translation which consist of only one or couple of words that occur in reference translations, as can be seen in the example below.

Intuitively, after a word from the reference translation has occurred it should not be considered in the calculation anymore. This intuition is formalized as the

modified unigram precision. It is computed in the following way: first, the maximum number of occurrences of a word in any reference translation is counted; then the total count of every candidate word is replaced by the maximum reference count, added up and divided by the initial total number of candidate words. As a result, the sentence which may receive high precision score will receive more realistic evaluation measured by modified precision score, as can be seen in the example below.

Candidate: of of of of of of of of of of

Reference: London is the capital of England and of the United Kingdom of Great Britain and Northern Ireland.

Precision: 1

Modified unigram precision: 3/10

Similarly is computed modified n -gram precision score for any n , but n -gram counts are collected instead.

XXX TODO: Sentence length

XXX TODO: The formula

1.4 Multi-target machine translation

1.4.1 Multi-lingual machine translation

With constant improvement of neural MT systems performance researchers started to experiment with incorporating multiple source and/or target languages into one model, and the results are promising: having $L1 \rightarrow L2$ and $L2 \rightarrow L3$ non-parallel corpora makes possible to train a model that can produce $L1 \rightarrow L3$ translation of decent quality; having high-resource $L1$ and low-resource $L2$ from the same language group helps increase Source \rightarrow $L2$ translation quality with pre-training on Source \rightarrow $L1$ data. In general, in some situations using more target and/or source languages in one translation model may not only insignificantly decrease its performance but also to improve it.

Even if the concept of combining multiple languages into one model and possible outcomes of such combination may seem intuitive, there exist multiple approaches of how exactly this might be performed. As for current time, Dabre et al. [2019] categorizes MNMT (multi-lingual neural machine translation) in the following way (Figure 1.2):

Multivay translation. The goal is constructing a single NMT system for one-to-many, many-to-one or many-to-many translation using parallel corpora for more than one language pair.

Low or Zero-Resource Translation. Large amount of parallel texts of high quality is available for most of European languages. However, it is not true for most of other languages in the world. Three main directions have been studied these cases. *Transfer learning*: Transferring translation knowledge from a high-resource language pair to improve the translation of a low-resource language pair. *Pivot translation*: Using a high-resource language (usually English) as a pivot to translate between a language pair. *Zero-shot translation*: Translating between language pairs without parallel corpora.

Multi-Source Translation. Important documents and internationally popular books have been translated into many languages. Having the source side represented by multiple languages may increase translation quality in general or help to remove ambiguities present in one or another source language (e.g. cases, noun genders, etc.).

1.4.2 Massively multi-lingual machine translation with complete sharing

In Aharoni et al. [2019] models with up to 103 languages were tested. English centric in-house dataset was used to train En→Any and Any→En multilingual models. The average number of examples per language pair is 940k: for 13 out of the 102 pairs there were less than one million examples available. All languages from 5-to-5 model are present in 25-to-25, same is true for all languages from 25-to-25 with respect to 50-to-50 and so forth. In all cases they trained large Transformer model with 473.7M parameters. As can be seen on Table 1.2, the quality of translation is significantly worse when model is trained to translate more languages. However, it is worth reminding that this many-to-many experiment may have different results due to many-to-one direction present in it.

The decrease of model’s performance with adding more target languages is clearly shown in Aharoni et al. [2019].

	En-Ar	En-Fr	En-Ru	En-Uk
5-to-5	12.42	37.3	24.86	16.48
25-to-25	11.77	36.79	23.24	17.17
50-to-50	11.65	35.83	21.95	15.32
75-to-75	10.69	34.35	20.7	14.59
103-to-103	10.25	34.42	19.9	13.89

Table 1.2: **BLEU scores for translation in one direction (part of Table 7 from [Aharoni et al., 2019])** . Model trained on 5-to-5 English centric dataset (English to any and any to English) scores 12.42 BLEU for English-Arabic test set. Every language from 5 languages of 5-to-5 data set is included into 25-to-25 set, as well as every language from 25-to-25 data set is included into 50-to-50 and so forth.

1.5 Conclusion

In this chapter we introduced theoretical and historical background for this work. Firstly, we took a short walk through the history of machine translation. Then we described the most used type of NMT models – self-attention *Transformer* model. After that we went over the history of translation evaluation in general and the most used method of automatic evaluation – BLEU – in particular. In the end multi-lingual neural machine translation was reviewed with more detailed view into ‘complete sharing’ scheme.

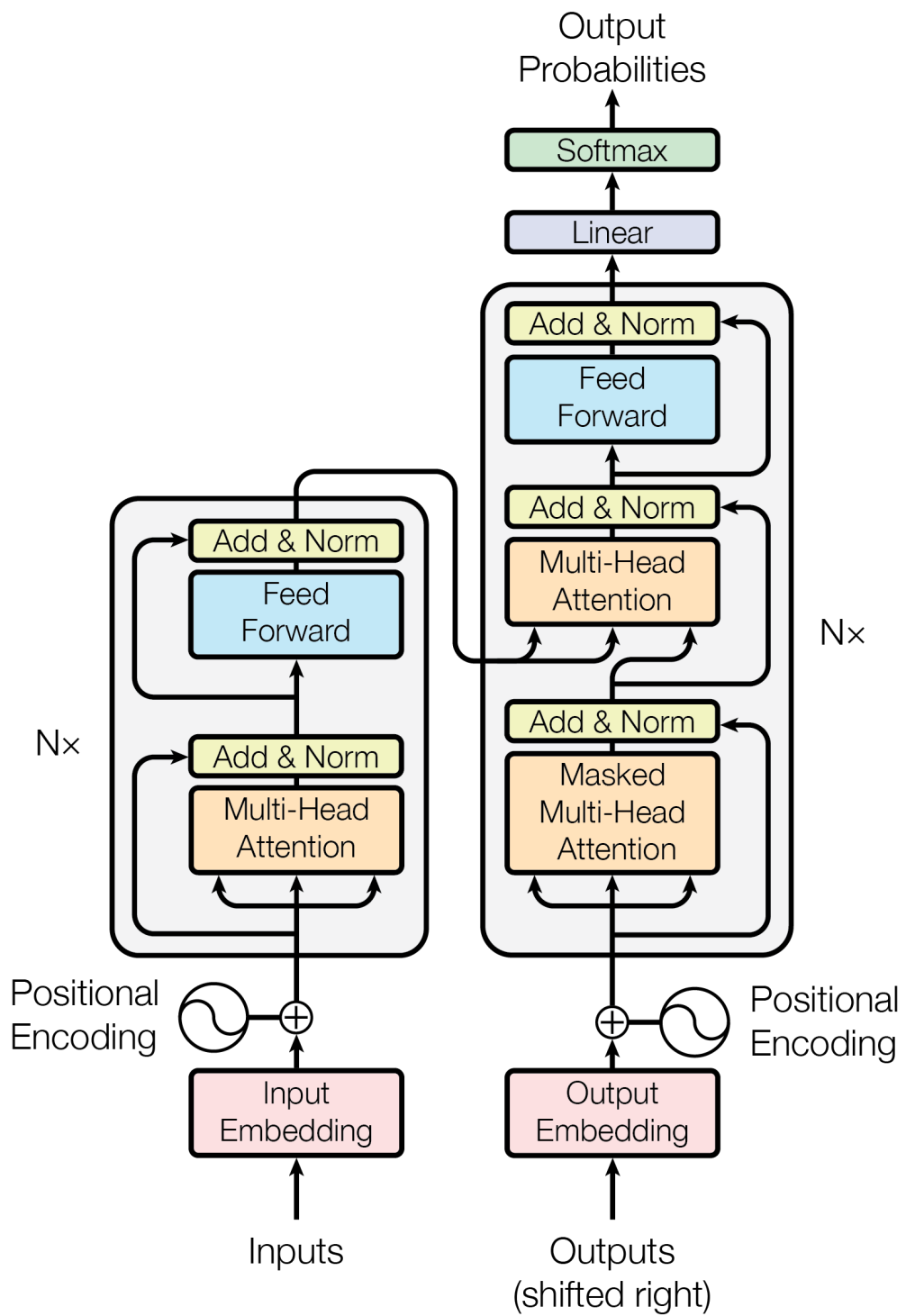


Figure 1.1: **Transformer model architecture.**

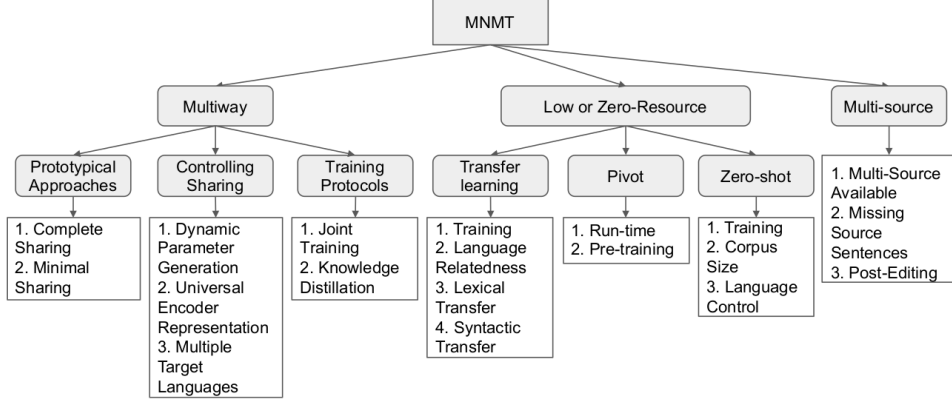


Figure 1.2: **MNMT research categorized.** According to resource scenarios and underlying modeling principles. By Dabre et al. [2019]

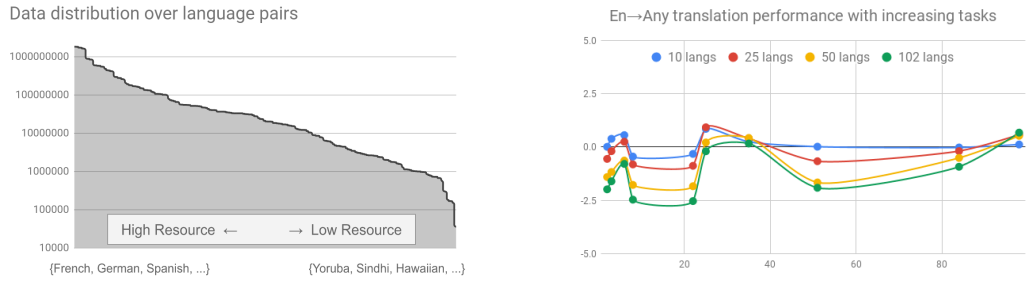


Figure 1.3: **Tranlsation performance for 102 languages from Arivazhagan et al. [2019]** . Axis X is shared between left and right plot. On axis X there are languages sorted by amount of training data. Left: amount of training data (axis Y) for a language. Right (best viewed in color): Effect of increasing the number of languages on the translation quality. On the axis X the languages are sorted the same way as on the left plot. The points visualized are 10 languages that are present in all setups from $En \leftrightarrow 10$ to $En \leftrightarrow 102$.

2. Experiment setup

In this chapter we describe the data used for experiments, training setup and experiments that were run to answer the questions asked in this thesis.

2.1 **XXX TODO: Questions and constraints**

Constraints:

Translation quality for multi-lingual system is insignificantly worse than for mono-lingual one-to-one translation system.

Maximum possible target languages are combined in one model.

Questions:

How *in average* adding one more randomly selected target language to the multitarget model affects its En→De performance?

How is it different if we add a linguistically similar, not randomly selected language?

How is adding one more language from the same language family or group *in average* affects translation performance for selected language pair (e.g. En→De)?

2.2 Experiments

2.2.1 Starting point

In Johnson et al. [2017] authors proposed a way to build a multi-lingual machine translation model without any changes to the *Transformer* architecture. In fact, the only change should be performed on the input data. To make the *Transformer* model process multi-lingual data the language tag is added to the source sentence. For example, the following En→Cz sentence pair:

Hello world! → Ahoj světe!

is modified to:

<2cs> Hello world! → Ahoj světe!

With given method it is possible to produce translations in multiple languages using the same model just by altering the prepended target language tag. It was also demonstrated that this method slightly improves translation quality for low resource languages when compared to monolingual translation model.

In this and the following (Arivazhagan et al. [2019], Aharoni et al. [2019]) papers from Google many different cases are tried and described. However, in each setting there is usually only one model of each kind considered. For example, when in [Aharoni et al., 2019] authors compare 5-to-5, 25-to-25, 50-to-50, etc. models, there is only one 5-to-5 model, one 25-to-25, etc.

Target language tags do not affect BLEU: Siddhant et al. [2019]
mNMT models en-to-4 and 4-to-en trained; 1) <2xx> added to the source;
2) target language encoded separately. BLEU scores are comparable using both approaches.

2.2.2 Proposed experiments

Bilingual baseline

It is important to have a reference point to be able to reason how does a change to the system affects its performance.

As a naïve and simple solution, so called baseline system, should be used a bilingual system with the same settings as the systems in further experiments.

Having a performance measurement for a bilingual system in some translation direction allows us to say about improvement or deterioration of results after implementing a change.

n-lingual baselines (random)

Multilingual models with random set of languages. The purpose is twofold: to show BLEU score decrease with increasing number of target languages and to serve as a baseline for multitarget models with target languages grouped by in non-random way, e.g. by language group or linguistic similarity.

Group by language group

If all target languages are from one language group we expect to observe better translation quality comparing to multilingual baseline results with randomly selected target languages. This is expected due to shared parts of vocabulary (todo: expand with examples) and linguistic properties (again, expand with examples). Germanic group: da, de, is, no, nl, sv. Slavic with cyrillic script: bg, mk, ru, uk. Slavic: bg, cs, hr, mk, pl, ru, sk, sl, sr, uk

2.3 Dataset(s)

2.3.1 TO EDIT: English to 36 languages

To observe effects of linguistic similarity of target languages it is important to examine enough possible variations of those. The OPUS dataset (Tiedemann [2012]) is an open and free collection of texts that covers more than 90 languages with data from several domains.¹

For our experiments the source language is English only.

Sampling and splitting of the data is the one used for ELITR project.² For each of language pairs and each sub-dataset the data was splitted to training, validation and testing sets. For each of the two latter sets 2000 random sentences were selected and the rest of the data remained for the training set. In cases where the sub-dataset contained less than 16000 sentence pairs no data went

¹Available at <http://opus.nlpl.eu/>

²https://elitr.eu/wp-content/uploads/2019/07/D11.FINAL_.pdf

to the validation set. Later for each language pair there were 1000000 sentence pairs sampled from all training sub-sets. Firstly, if available, the sentences were taken from Europarl, then EUbooks, OpenSubtitles, and then all remaining sub-datasets. The same procedure was used to sample x000 of validation set sentences per each language pair. The test sets were left separate, so that the result on each domain would be observable.

Later an overlap in the source side of different language pairs was found. Although this would not directly lead to unfair increase of the test score, such sentence pairs were removed from the training sets. This filtering decreased the amount of sentence pairs to 0.85-0.95 millions per language pair.

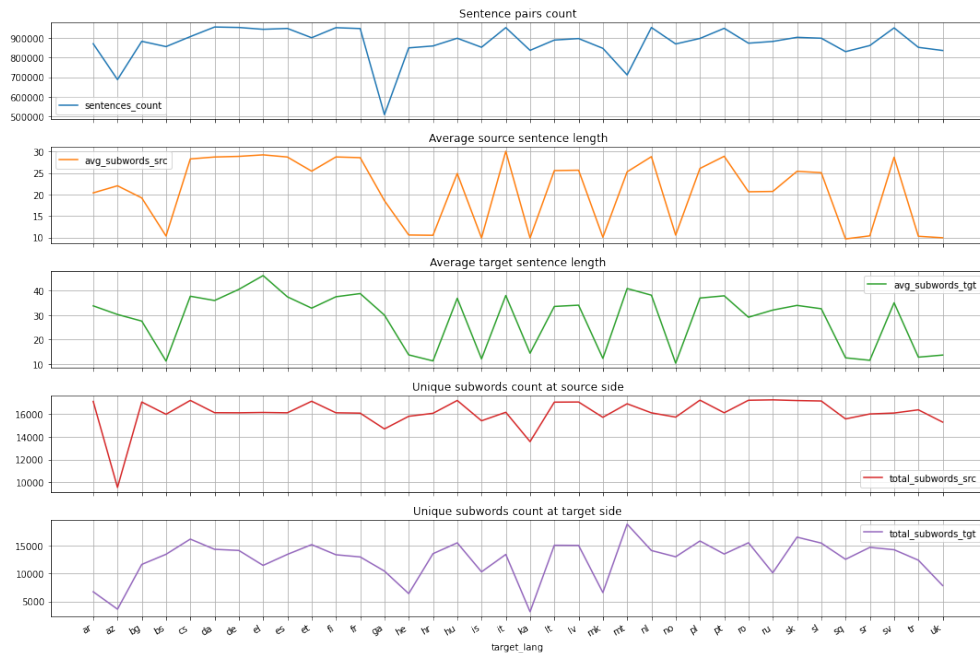


Figure 2.1: **Training data language statistics.** Languages are on the X axis sorted as in appendix. From top to bottom: total number of sentence pairs in training set per language, average amount of subwords per sentence on the source side, the same on the target side, total amount of unique subwords for this target language on the source side, the same on the target side.

2.3.2 XXX TODO: UN parallel corpus: English to 5 languages

XXX TODO: as I show en-to-5 results in RANDOM section, this DS should be here Eisele and Chen [2010]

2.4 Method

In this section we describe how the models are trained, which metrics are collected and how are they analyzed.

group	subdataset names	description
1	Europarl/vx, DGT, MultiUN, EUbookshop, JRC-Acquis, ECB, EMEA	Proceedings and documents from Europarl, UN, etc.
2	NewsCommentary, GlobalVoices, WMT-News	News articles and commentaries
3	OpenSubtitles, Tatoeba	Short sentences, human speech, general domain
4	OpenOffice, PHP, KDE4, Gnome	Software documentation or interface elements
5	Bible, Tanzil, Books	Other: completely different domain, dated vocabulary automatically aligned sentences, etc.

Table 2.1: **Groups of subdatasets in OPUS.**

2.4.1 TO EDIT: Training

For example, let us take $\text{En} \rightarrow \{\text{Fr}, \text{De}\}$ setup, which means that the model to be trained should take a source sentence in English and produce translation either in French or in German. The language of model’s translation depends on the target tag at the beginning of the input sentence, i.e. `<2fr>` tag in source sentence leads to French translation.

To train such a model, only related sentence pairs are subsampled from the whole training set. In this case, from the whole training set we select only those sentence pairs which source side starts with tags `<2fr>` or `<2de>`. Such subsampled dataset is then used to train the model.

During the training procedure, once per specified number of updates occurs the checkpointing of the model. The model weights are saved to the disc and

number of metrics are logged.

Those metrics are:

- training loss value (mean value for all updates since last checkpoint)
- current learning rate value
- training speed (processed words per second)
- training time since last checkpoint
- current number of updates

Hardware usage should also be recorded if possible:

- GPU usage
- CPU usage
- memory usage
- disc I/O
- network I/O

The hardware metrics are not important for model's evaluation but may help early spot mistakes like underuse of GPU or CPU, lack of RAM, etc.

2.4.2 TO EDIT: Validation

The validation set is used to track model's performance on an unseen set of data and to perform early stopping. Once per specified number of steps the validation occurs: validation metrics are recorded, for any metric which value was improved current model weights are saved as best model by this metric. If early stopping condition occurred then the training process is stopped.

For any model the validation set is constructed of the big validation set by selection only relevant sentence pairs in the same way, as the training set.

For the example setup from above, $\text{En} \rightarrow \{\text{De}, \text{Fr}\}$, the validation set consists of equal amount of $\text{En} \rightarrow \text{De}$ and $\text{En} \rightarrow \text{Fr}$ sentence pairs. E.g. if in the complete validation set there are 1000 sentence pairs for each of possible target languages, then for $\text{En} \rightarrow \{\text{De}, \text{Fr}\}$ model the validation set will contain 2000 sentence pairs, and for $\text{En} \rightarrow \{\text{De}, \text{Es}, \text{Fr}\}$ it will contain 3000 sentence pairs.

For the validation set we collect not only the loss function value but also the metric of interest, which is BLEU score. However, the BLEU of the whole model's validation set is not something we are interested in. For the discussed example we collect `bleu:fr` and `bleu:de` scores which represent BLEU scores for French and German parts of validation set. E.g., to compute `bleu:fr` we select only $\text{En} \rightarrow \text{Fr}$ sentence pairs from the validation set.

Also, an aggregated value of the `bleu:xx` scores is also recorded and may be used for early stopping: ending the training process when the metric is not improved during last N validation steps.

All together, the following validation metrics are recorded after the validation step:

- loss function value
- bleu:xx which is BLEU score for each of model's target languages
- aggregated value of all bleu:xx values
- translation time of the model's validation set

2.4.3 XXX TODO: Finishing the training

When should we stop the training? It is not possible to say precisely when did the model acquire its best performance because of stochastic nature of the training algorithm (SGD). Because of that we need to use some method how to decide when training process should be stopped.

Number of epochs

XXX TODO: First occurrence of a term to be italic The easiest approach is to specify number of epochs after which the training is stopped. This could be a good solution for the case when all models that will be compared are trained on the same amount of data from the same domain. But in our case, adding one more target language adds a constant amount of sentence pairs to the training set. Roughly, if the number of epochs is specified as a stop condition, a bilingual En→De model will see the German training data x times, when multilingual En→De, Fr, Es will only see the German training data $x/3$ times.

Early stopping

Early stopping is a regularization technique used to avoid possible overfitting of a model on the training data. In general, it works in the following way: after every validation step it checks if the metric value improved during last N validations. The metric to be controlled and number of validation steps N are the parameters of this method (Figure 2.2).

Another situation is even more probable: model's validation performance is either stalled or slightly improved (Figure 2.3). In this case early stopping helps to avoid unnecessary spendings on computational resources.

In our case we could use early stopping to ensure more equal conditions for models with different sizes of training data. Suitable number N could be found experimentally, but which metric should be used?

XXX TODO: cross-entropy, perplexity; they may not represent the model's performance, what about BLEU

Given that the task is to train a model that is as good as possible in **every** of its target directions, the BLEU score of the whole translated validation set for this set of languages does not say anything about the model's performance in each specified translation direction.

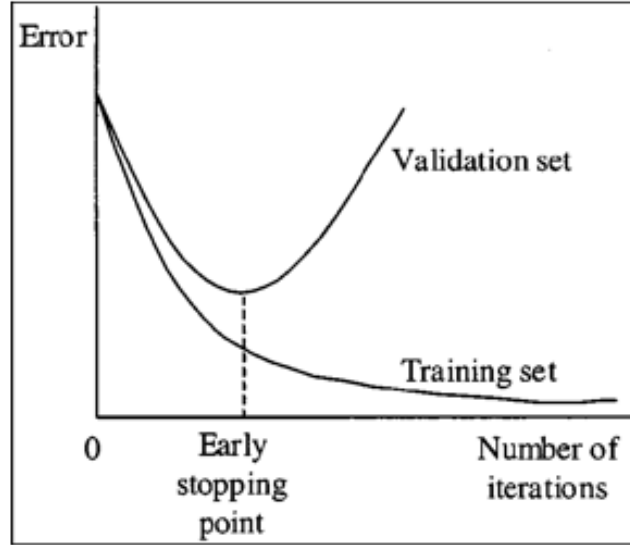


Figure 2.2: **Early stopping to prevent overfitting.** At the 'early stopping' point the model's performance on unseen validation set of data does not improve anymore. Further training leads to poorer performance on unseen data. Stopping the training at this point results in better model's performance on unseen data.

Aggregated value of BLEU scores

Therefore, we should use separate BLEU scores which represent model's performance in each of translation directions. The most intuitive and naive way is to compare BLEU scores for each target language. But if we will try to set up an early stopping on each of the values, the following may happen: BLEU for one of the directions will stay the same for long enough to trigger early stopping, but for another direction the BLEU score is still growing.

For example, suppose the model is being trained for $\text{En} \rightarrow \{\text{De}, \text{Fr}\}$ and validation happens after each 1000 update steps. At the step 20000 for some selected German validation set the model's $\text{En} \rightarrow \text{De}$ performance measured by BLEU score equals 20, for the $\text{En} \rightarrow \text{Fr}$ direction its BLEU score is 15. After 10 more validations, model's performance for $\text{En} \rightarrow \text{De}$ is still equal to 20 BLEU and did not grow, but for $\text{En} \rightarrow \text{Fr}$ it has already reached 25 and is still increasing. Early stopping will be triggered, and even though the performance in $\text{En} \rightarrow \text{De}$ direction might have reached its possible maximum (assuming that the early stopping period of 10 checks was selected efficiently) its $\text{En} \rightarrow \text{Fr}$ performance might still improve.

XXX TODO: Geometric mean: good for multiple metrics with different scale

2.4.4 TO EDIT: Testing

After the training is finished, the received models should be evaluated on unseen test data. For experiments with en-to-5 dataset (Section 2.3.2) the test sets are created in the same way as validation sets. For en-to-36 dataset (Section 2.3.1) the test set is divided on subsets by the source dataset. It means, that each of source datasets (like OpenSubtitles/v11, Europarl/v7, etc.) there exists a separate test set. So after translating the test set for each of target directions of the model the following record is created:

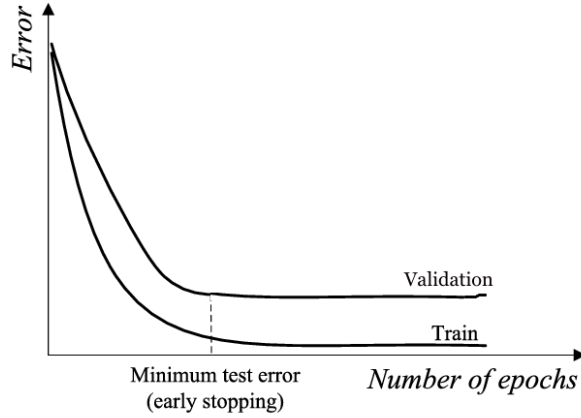


Figure 2.3: **Early stopping as the model is not improving.** Even though the metric value on the training set is still slowly improving, the its value on the unseen validation set is stalled. Further spending of computational resources is unjustified.

- model name
- source language
- target languages
- tested target language
- BLEU score for this part of translation
- metric, based on which the best model was saved
- dataset name (for en-to-36)

Let us return to the example setup is $\text{En} \rightarrow \{\text{De}, \text{Fr}\}$ and suppose the reported validation metrics are the mean loss function value on test set and 'translation' (geometric mean of all reported BLEU scores, see Section 2.4.2). After the training is finished, there will be two models: best by loss value and best by 'translation'. For each of those two records are created: for $\text{En} \rightarrow \text{Fr}$ translation and for $\text{En} \rightarrow \text{Es}$. In total, 4 results are recorded.

If the model was trained and tested on en-to-36 dataset, than 4 times n records are created, where n is number of OPUS subdatasets from which the data was sampled.

2.4.5 **XXX TODO: Analysis**

After required set of models is trained and their test metrics are collected, data should be analysed.

2.5 Training tools

In the following section we describe the tools that are uset to implement what was shown in Section 2.4.

2.5.1 Tool kits

There exists a number of different tools that can be used for training a NMT model. General purpose deep learning programming libraries like Tensorflow³ and PyTorch⁴ are most popular for deep learning related research. With their help it is possible to construct any of today’s state-of-the-art NMT models; pre-built and pre-trained models are initially present in such frameworks, but it is also possible to describe a model from scratch.

Another option is presented by specialized NMT tool kits. They usually contain efficient and tested implementations of NMT models as well as some of usefull preprocessing tools. For the experiments described in 2.2 there is a need to train significant amount of models with the same architecture and settings but different datasets. Due to that fact, in this work the use of specialized NMT tool kit is more suitable. Let us consider the foolowing list of broadly used tool kits as for year 2020, presented in Koehn [2020]:

- OpenNMT (based on Torch/pyTorch)⁵
- Sockeye (based on MXNet)⁶
- Fairseq (based on pyTorch)⁷
- Marian (stand-alone implementation in C++)⁸
- Google’s Transformer (based on Tensorflow)⁹
- Tensor2Tensor (based on Tensorflow)¹⁰

We chose *MARIAN-NMT* tool kit¹¹ as a fast solution with stable and efficient *Transformer* Vaswani et al. [2017] implementation, minimum of third-party dependencies, and ability to train models on multiple GPU units in parallel.

2.5.2 Computational cluster

In the experiments proposed above the expected number of models to be trained is quite big. First of all, there should be 36 models for *mono-target baseline* for En→36 dataset. For *multi-target random* experiment the number is much bigger. For example, let us consider a case with En→3 models - each model translates from English to 3 target languages. Specifying that each of 36 target languages from En→36 dataset should appear at least in 3 En→3 models, series of random generation of En→3 setups gave the smallest amount of such setups equal to 44. For En→5 case with 5 target languages in each model and with the

³<https://tensorflow.org/>

⁴<https://pytorch.org/>

⁵<https://opennmt.net>

⁶<https://github.com/aws-labs/sockeye>

⁷<https://github.com/pytorch/fairseq>

⁸marian-nmt.github.io

⁹<https://github.com/tensorflow/models/tree/master/official/transformer>

¹⁰<https://github.com/tensorflow/tensor2tensor>

¹¹Junczys-Dowmunt et al. [2018]

same restriction of minimum occurrence the same procedure gave the minimum amount of needed models equal to 34.

To be able to train large number of models in a reasonable amount of time we needed to use computational cluster with GPU cards. The computational clusters available at the institution are operating under SGE¹² scheduling software and are equipped with GPU cards with minimum CUDA *compute capability* 6.1.

Considering data storage quota limitation and high utilization of computational resources by the cluster’s users, the following training pipeline was designed:

- Prepare task list
- Iterate over the list working with at most N tasks in parallel
- For each task
 - Subsample the dataset taking only those sentence pairs with target languages specified in the task
 - Run the training procedure for limited amount of time (e.g. for one hour only) starting with previous checkpoint if it already exists
 - Regularly compute metrics on the development set and report them
 - On the event of evaluation on the development set save the best model for each metric
 - After time is out the training is stopped and subsampled datasets are removed
- If for next selected task the model is already trained then select next task from the list
- If for next selected task the model is currently being trained then decrease the number N of tasks processed in parallel

2.5.3 Inspecting the training process

As the number of models trained and being trained is growing, monitoring of the training process becomes more and more complicated. If the experiments are also being run on different computational clusters it becomes very possible that a parameter mistakenly set up to different value or a corrupted dataset, or even hardware version may lead to an unexpected difference in results.

To address these and other issues that may occur during the training process we use Weights&Biases¹³ experiment tracking tool. Its main features that are useful in this prospective are following:

- Metric visualization
 - Training and validation loss curves (Figure 2.4 left subplot)
 - Scatter plots (Figures 2.5 and 2.4 middle subplot)

¹²<https://arc.liv.ac.uk/trac/SGE>

¹³Biewald [2020]

- Artifact storage
 - Model checkpoints storage
 - * stores 'heavy' model files which cannot be stored in *git*
 - * along with *git* it makes possible to move training to the different computational cluster system
 - Sample translations of validation set
 - * helps to observe improvements of translation quality in time
 - * lets verify that model is actually produces meaningful translation
- Customizable reports
- Hardware utilization

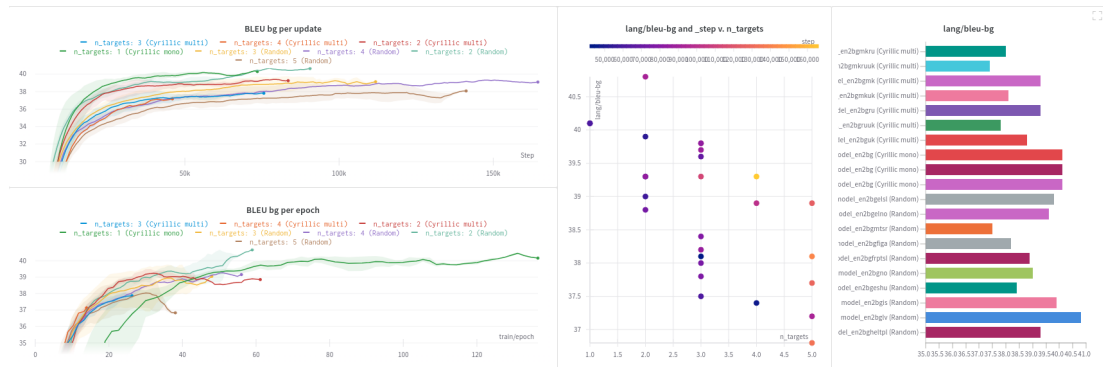


Figure 2.4: **Single language results visualization: models other target languages are randomly selected vs. those selected from similar group of languages.** XXX TODO: split it in print-friendly way Here is a part of interactive report for 'slavic languages with cyrillic script vs. random' experiment. In this specific case models' performance on Bulgarian part of validation set is compared.

Left: BLEU score for En→Bg translation direction is monitored with training step on X axis (top) and training epoch (bottom). Each curve represents mean value (line) and its min/max value (range) at the point of time of multiple models' results. Models are grouped by number of target languages and experiment subgroup (En→Bg, En→Slavic and En→Random).

Middle: Number of targets (axis X) vs. BLEU on En→Bg validation set (axis Y) vs. update steps (colos with scale at the top).

Right: Individual models' En→Bg validation BLEU scores.

2.5.4 XXX TODO: Model settings

The initial parameter selection is made with respect to Popel and Bojar [2018]. First of all, the hyperparameters of MT model are tuned on couple of language pairs from one dataset. The parameters leading to the same result in shorter time were preferred. Then the selected parameters were used on all experiments with the dataset.

Tuning early stopping on early runs

The initial early stopping setting was that after 5 consecutive validation steps without improvement of validation loss value the training process stopped. However, during the training of first couple of bilingual models the following situation have happened quite often. Further improving performance on validation set by couple of tenths of BLEU points took as much time as reaching the optimal state.

On the Figure 2.6 can be seen that path from the beginning of training to optimal point B (26.9 BLEU) took as much time as its further improvement by 0.2 BLEU at point D (27.1 BLEU). However, there were certain models with a bit bigger improvement after a much longer time, e.g. 0.8 BLEU points on Figure 2.7.

Possible situations of this kind were discussed in Section 2.4.3. After considering also some of preliminary multilingual runs the 'patience' parameter of early stopping was set to 15. After 15 consequent validation steps without metric improvement the training process is stopped.

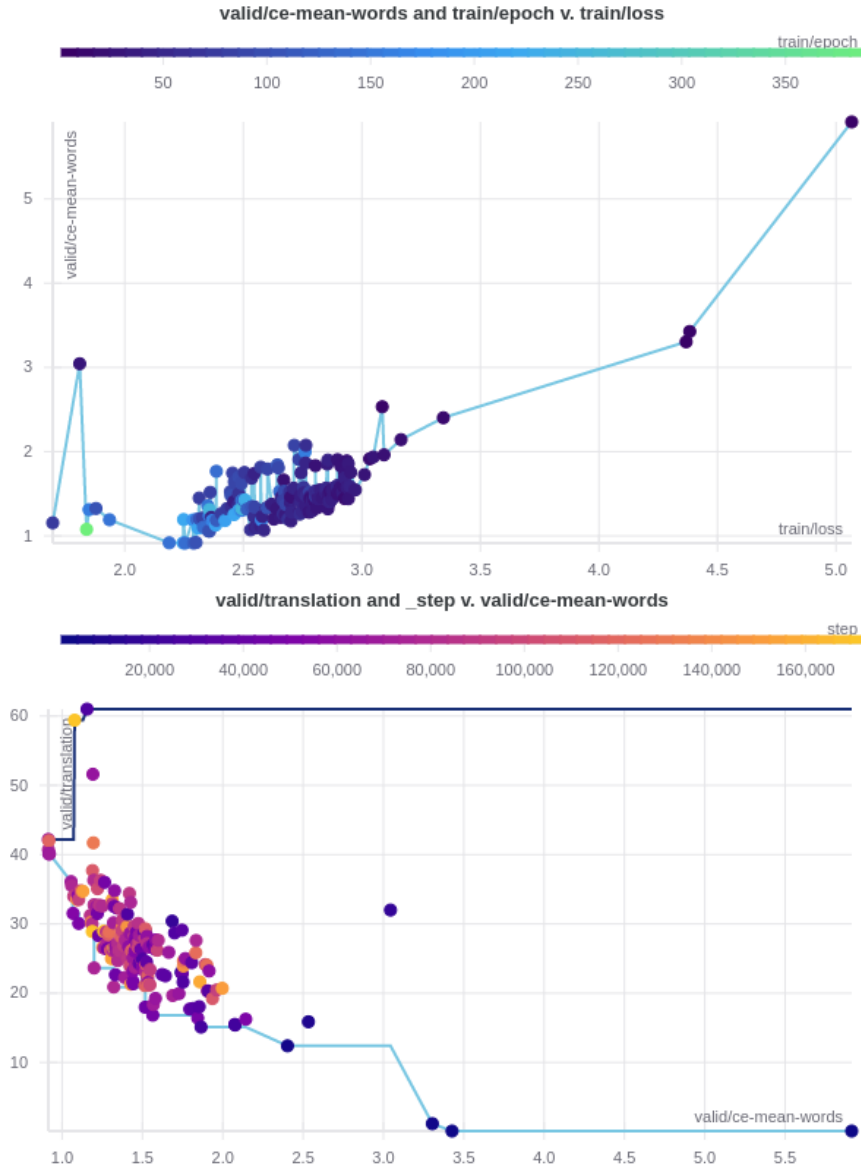


Figure 2.5: **Overall model convergence dashboard.** On these two interactive graphs each point represents one model. Models that are currently training are visualized here together with completely converged models and those which training process is currently on hold.

Left: axis X represents the training loss value, axis Y represents the value for the same loss function calculated on the validation set. The color of each point represents current training epoch for the model. Normally for any model the point moves from top right part of this graph to the bottom left part, representing both training and validation loss being gradually decreased during the training procedure. The point that moves to the middle left part of the graph may signalize about either overfitting of the model on training set, or difference in data distribution in training and validation set, or else some mistake in training settings.

Right: on this plot loss value on the validation set (axis X) is compared with geometric mean of BLEU scores for each of target languages. For any model during the training procedure its point usually moves from bottom right corner into the cluster of other points. Model which point 'arrives' to any other location than the cluster may need special attention.

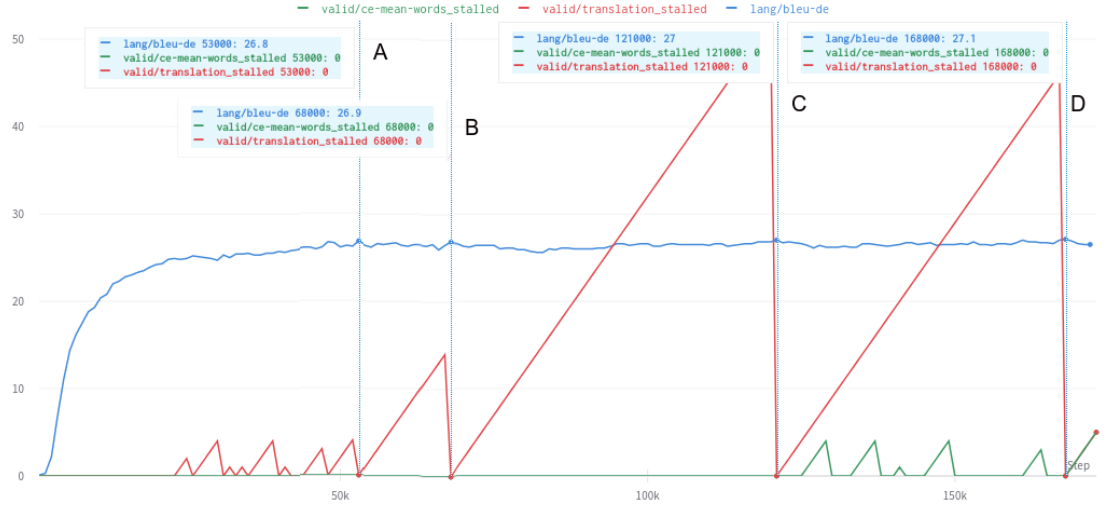


Figure 2.6: **Example change of model's performance on validation set in time.** Preliminary En→De model.

Blue: validation metric (value on the left axis in BLEU)

Red: validation metric (BLEU) stalled. Each consecutive validation step when the metric is not improved this value is incremented by 1. When the metric is improved this value is reset to 0.

Green: loss function value on validation set is stalled. Same logic as for *Red*.

BLEU score values at the points of improvement: *A* – 26.8, *B* – 26.9, *C* – 27.0, *D* – 27.1.



Figure 2.7: **Small improvement during long training.** In this case (En→Fi) the difference is a bit more visible: 21.3 at the first point and 21.9 at the best. Colors and scales are the same as at Figure 2.6.

3. Bilingual and multi-lingual baselines

In this chapter we describe the baseline experiments. Bilingual baselines are needed to specify the starting point: how good can model perform on specific translation direction for each test set.

After bilingual results are collected and inspected, it is time for multi-lingual baselines. As multi-lingual baselines we consider models with randomly selected set of target languages. This way we can see how much adding more target languages to the model changes its performance on the same specific translation direction.

Most of experiments are done on en-to-36 dataset with couple of additional experiments on en-to-5.

3.1 XXX TODO: Bilingual baseline

We trained bilingual models on en-to-36 dataset and received number of values for each translation direction. Test results for relevant target directions (i.e. languages from 'Germanic' and 'Slavic with cyrillic script' from Section 2.2.2) are shown in Table A.1.

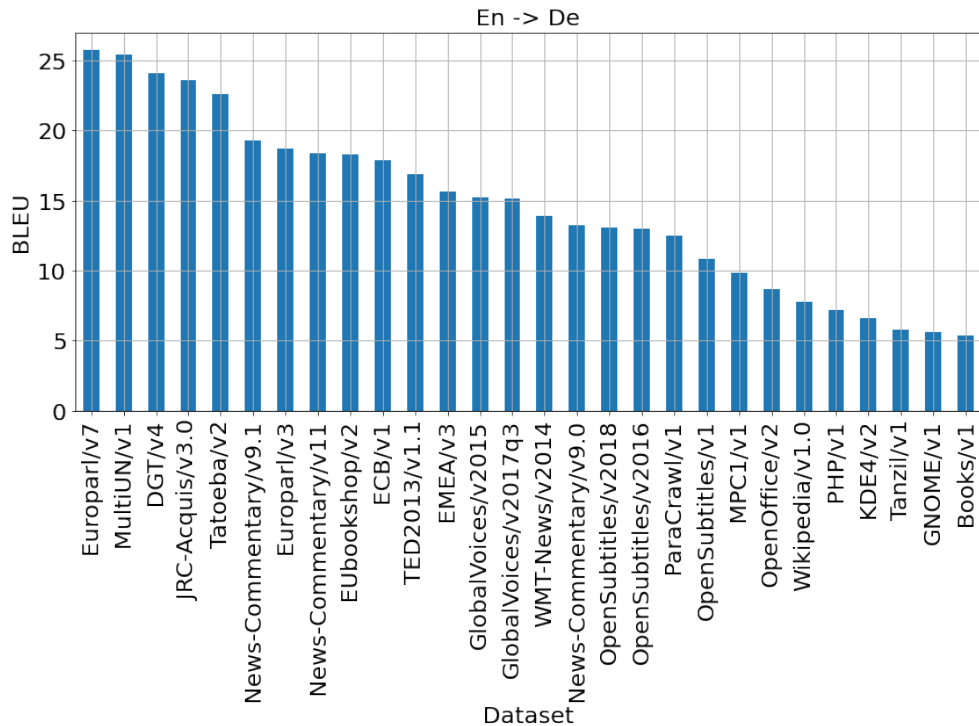


Figure 3.1: **En→De bilingual results.** Datasets on the X axis are sorted by declining BLEU score.

3.2 XXX TODO: Multilingual baseline

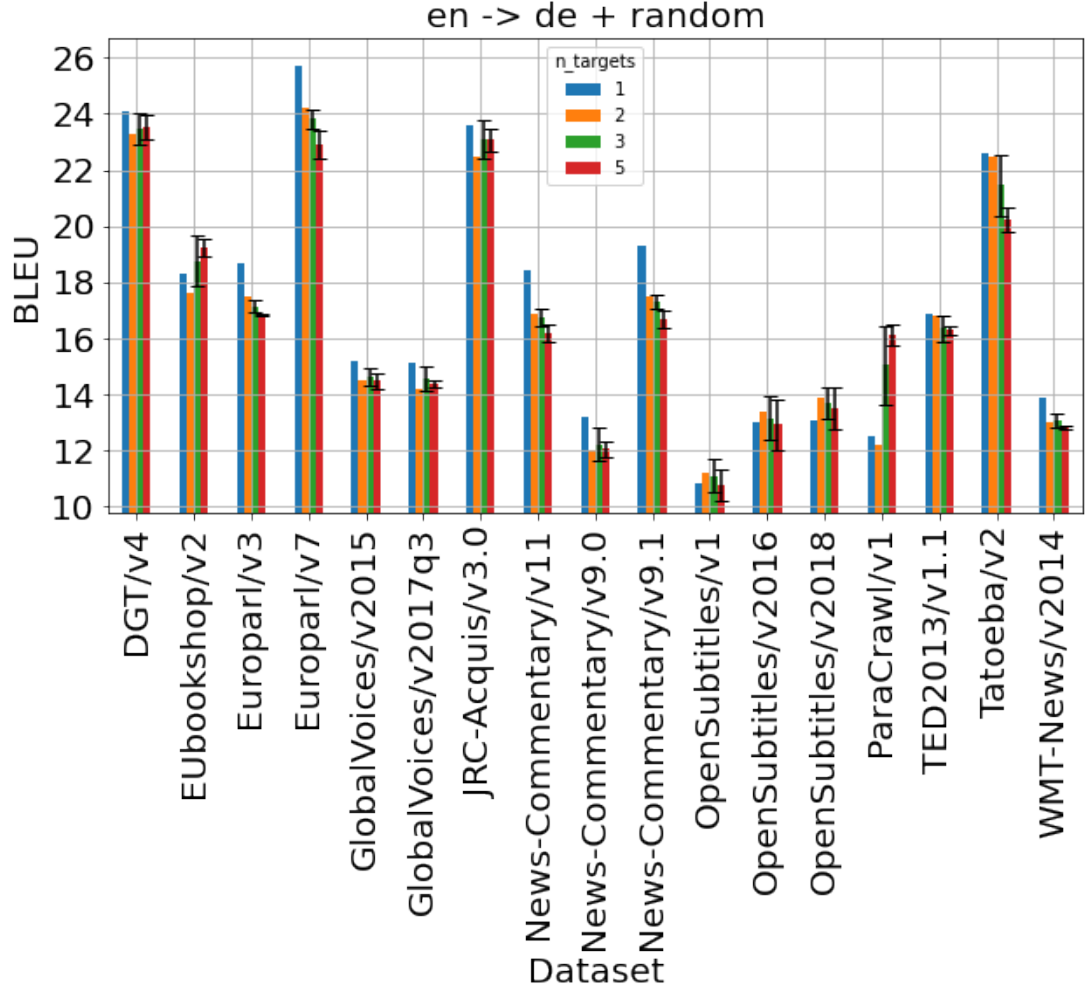


Figure 3.2: **En→De monolingual baseline results (RANDOM).** Datasets with BLEU lower than 10 are removed

E.g. as a bilingual baseline model for German target language the En→De model is trained. En→De,Fr, En→De,Az, En→De,Bg En→Bg, Az models provide us with 3 results for En→De direction 2-target baseline, as well as one value for En→Fr, two values for En→Bg and two for En→Az. These aggregated En→De, X results will be later compared with aggregated En→De, X1, X2 for 3 target languages, En→De, X1, X2, X3 for 4 target languages, where X1, ... Xi are randomly selected languages. Also in the next chapters these results will serve as a baseline for aggregated En→De, Y1, Y2... results of N target languages for languages Yi being from some group of languages similar to De.

3.3 Expected results

As we have seen in section 1.4.2, models with more languages in the mix usually perform slightly or significantly worse than bilingual ones. Also, for bilingual

baselines no significant change in performance is expected with adding the target language tags.

However, there might be different unexpected effects due to slight domain-wise differences in corpora content for different target languages.

3.4 Performance drop on massively multilingual setup

1-to-3, 5, 7, etc. models on en-to-36 dataset (0.9 mil. sentences per target language)

When the size of the model is fixed, adding more translation directions usually causes worsening of its performance. Multiple studies have shown this to be truth for many-to-many setup.

n_targets	mean	std	count	n_targets	mean	count	std
1	41.40	—	1	1	19.50	1	—
2	40.60	0.20	3	2	18.88	4	0.39
3	39.39	0.62	8	3	17.45	4	0.52
4	39.40	0.71	2	4	17.80	2	0.42
5	38.45	0.52	6				

(a) En→Bg for *Europarl/v7* dataset.

(b) En→Ru for *OpenSubtitles/v2016* dataset.

Table 3.1: **BLEU score change with adding target languages.** (a) First row: for mono-lingual En→Bg model test BLEU score is 41.40. Second row: for 3 (column *count*) En→Any models with two target languages (column *n_targets*) one of which is Bulgarian the mean BLEU score is 40.60 with standard deviation 0.20. (b): same way as (a)

3.5 Performance decrease on richer data sets

1 to 3, 4, 5 on UN corpus (much more sentence pairs per target language)

model	BLEU
m.esfrru t.ru	40.35
m.esru t.ru	42.16
m.frru t.ru	41.95
m.ru t.ru	44.20
m.esfrru t.fr	45.03
m.esfr t.fr	46.84
m.frru t.fr	45.66
m.fr t.fr	48.64
m.esfrru t.es	56.33
m.esfr t.es	57.94
m.esru t.es	57.31
m.es t.es	59.94

4. Group by language groups

4.1 Language groups

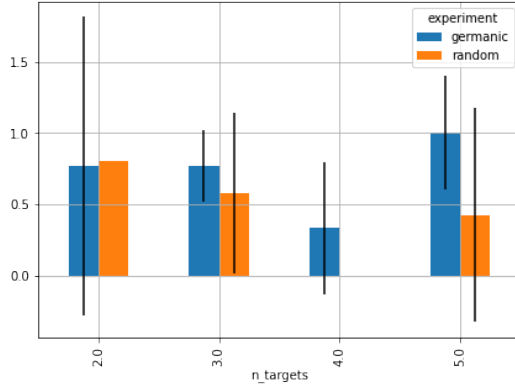
1 to 2, 3, 4, 5, etc. models on en-to-36 dataset (0.9 mil. sentences per target language) compared with random runs

4.1.1 Germanic group

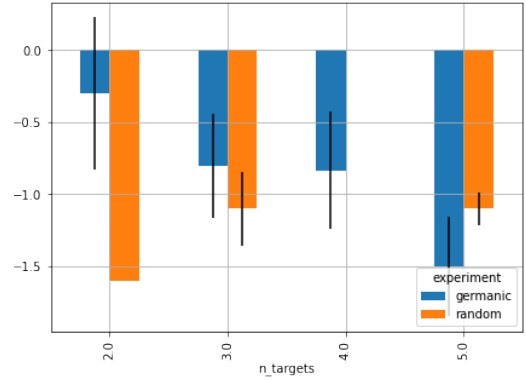
Here Germanic group consists of German, Dutch, Swedish, Danish, Norwegian and Icelandic. Models En→Germanic are compared to En→non-Germanic, where non-Germanic consists of any language except from the Germanic group. On Figures 4.1 and 4.2 some selected results are visualized along with vocabulary changes. Results for OpenSubtitles/v2018 mean the BLEU score on test set part sampled from OpenSubtitles/v2018. On both figures the subfigure (a) shows the result on spontaneous or pseudo-spontaneous speech transcripts, sub-figure (b) shows the result for prepared speeches or documents from Europarl or UN meetings.

In this case observations are twofold:

- For test sets with lower bilingual BLEU score adding more target languages to the model improves the score; adding related target languages improves it even more
- Adding more target languages improves translation result on test sets from spontaneous speech domain but worsenes it for prepared speech or documents.

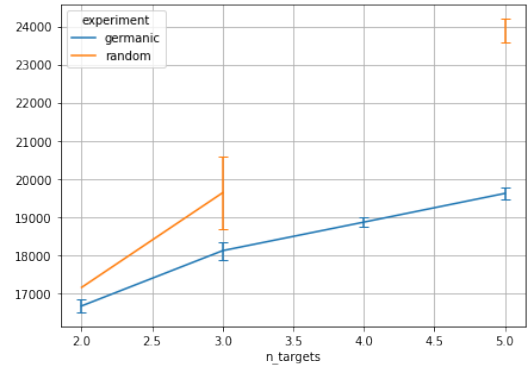


(a) OpenSubtitles/v2018, bilingual score: 13.1 BLEU

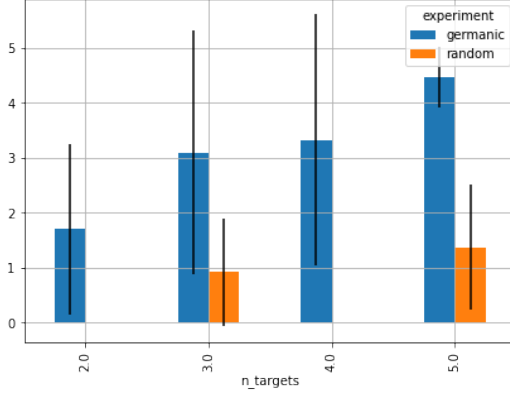


(b) MultiUn, bilingual score: 25.4 BLEU

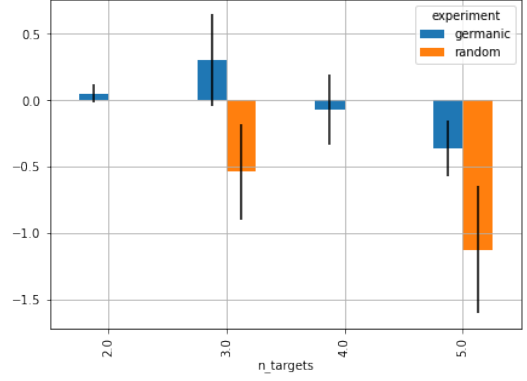
Figure 4.1: En→De BLEU score difference: Random vs. Germanic. On X axis - number of target languages. On Y axis - difference score comparing with monolingual BLEU. Black vertical lines show standard deviation. (a) Adding random target language as well as a related one slightly improves German translation score on speech transcript. (b) Adding neither a random nor a related target language helps with prepared speeches transcripts and documents in German. (c) Adding a related target language into the mix introduces less new unique subwords.



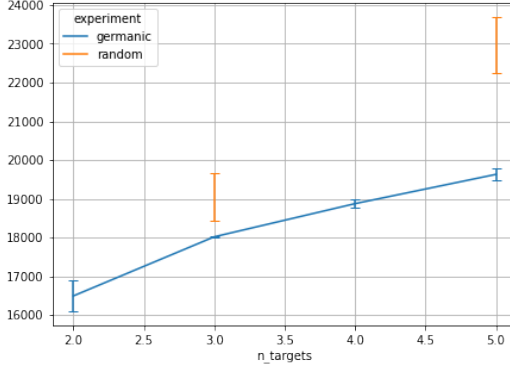
(c) Subword dictionary size used for target side



(a) OpenSubtitles/v2018, bilingual score: 15.6 BLEU



(b) Europarl/v3, bilingual score: 24.6 BLEU



(c) Subword dictionary size used for target side

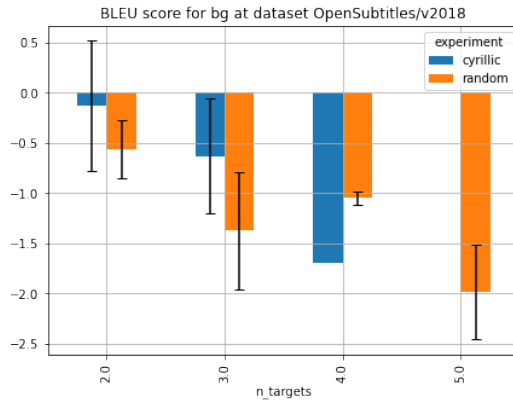
Figure 4.2: **En→Da BLEU score difference: Random vs. Germanic.** Axis are same as above.

(a) For OpenSubtitles test set which consists of human speech transcripts adding similar target language to the mix significantly improves the result. (b) For Europarl/v3 which consists of prepared speeches transcripts and documents adding more germanic languages to the mix did not worsened Danish translation quality unlike the case with German. (c) Adding random target language to the mix adds more subwords to the target subword dictionary

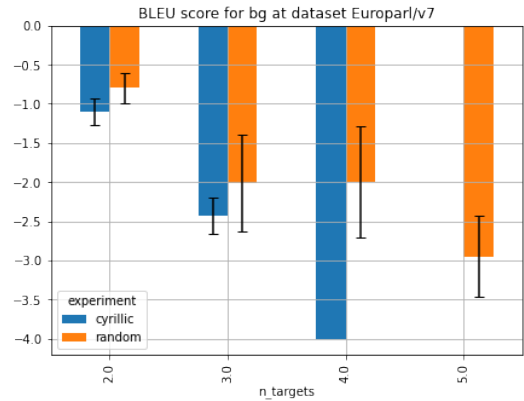
4.1.2 Slavic with cyrillic script

Here *Slavic with cyrillic script* group consists of Bulgarian, Macedonian, Russian and Ukrainian. Models En→Cyrillic are compared to En→non-Cyrillic, where non-Cyrillic consists of any language except from those from the group above. On Figures 4.3 and 4.4 some selected results are visualized along with vocabulary changes. Test sets for subfigures (a) and (b) selected the same way as in 4.1.1.

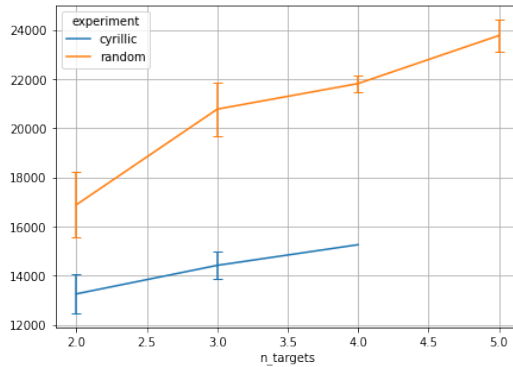
From the two opposite observations of 4.1.1 in this case the second one is observed: low results are getting slightly better, good results are getting slightly or significantly worse.



(a) OpenSubtitles/v2018, bilingual score: 23.7 BLEU

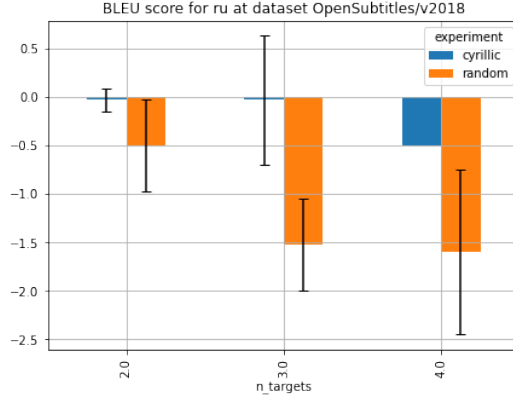


(b) Europarl/v3, bilingual score: 41.4 BLEU

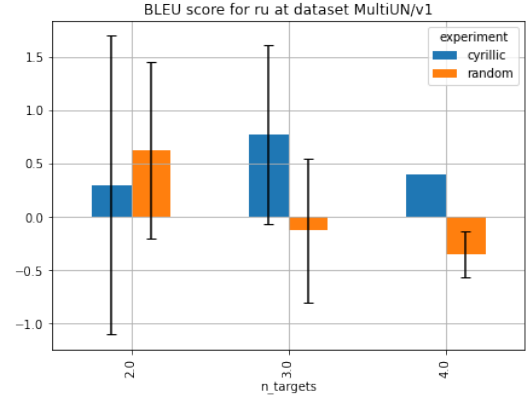


(c) Subword dictionary size used for target side

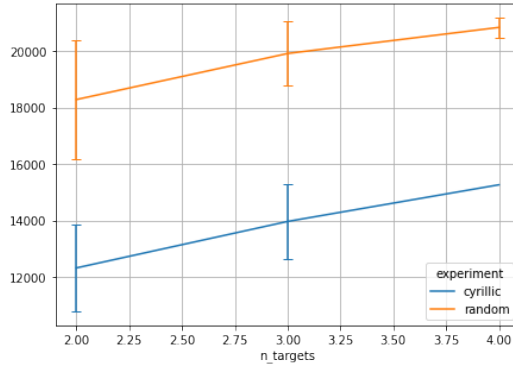
Figure 4.3: En→Bg BLEU score difference: Random vs. Slavic with cyrillic script. Axis are same as for Figures 4.2 and 4.1. There is not any data for Cyrillic and 5 targets as there are only 4 such languages in the en-to-36 dataset. Both (a) and (b) show significant decrease in translation quality. On (c) it is clearly visible how adding a random language with non-cyrillic script increases target subword vocabulary size.



(a) OpenSubtitles/v2018, bilingual score: 19.2 BLEU



(b) MultiUN, bilingual score: 14.6 BLEU



(c) Subword dictionary size used for target side

Figure 4.4: **En→Ru BLEU score difference: Random vs. Slavic with cyrillic script.** Axis are same as for Figures 4.2, 4.1 and 4.3. There is not any data for Cyrillic and 5 targets as there are only 4 such languages in the en-to-36 dataset. Both (a) and (b) show significant decrease in translation quality. On (c) it is clearly visible how adding a random language with non-cyrillic script increases target subword vocabulary size.

5. Discussion

5.1 Results

More languages in the mix: - share word ordering patterns - share vocabulary

More training data illustrates the properties of distribution better (some words are rare, some are often used)

5.2 Further work

Conclusion

Bibliography

- Roei Aharoni, Melvin Johnson, and Orhan Firat. Massively Multilingual Neural Machine Translation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3874–3884, Minneapolis, Minnesota, 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1388. URL <https://www.aclweb.org/anthology/N19-1388>.
- ALPAC. *Language and Machines*. National Academies Press, Washington, D.C., jan 1966. ISBN 978-0-309-57056-5. doi: 10.17226/9547. URL <http://www.nap.edu/catalog/9547>.
- Naveen Arivazhagan, Ankur Bapna, Orhan Firat, Dmitry Lepikhin, Melvin Johnson, Maxim Krikun, Mia Xu Chen, Yuan Cao, George Foster, Colin Cherry, Wolfgang Macherey, Zhifeng Chen, and Yonghui Wu. Massively Multilingual Neural Machine Translation in the Wild: Findings and Challenges. jul 2019. URL <http://arxiv.org/abs/1907.05019>.
- Loïc Barrault, Ondřej Bojar, Marta R Costa-jussà, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, Mathias Müller, Santanu Pal, Matt Post, and Marcos Zampieri. Findings of the 2019 Conference on Machine Translation (WMT19). In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 1–61, Stroudsburg, PA, USA, 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-5301. URL <http://www.aclweb.org/anthology/W19-5301><https://www.aclweb.org/anthology/W19-5301>.
- Lukas Biewald. Experiment tracking with weights and biases, 2020. URL <https://www.wandb.com/>. Software available from wandb.com.
- Ondřej Bojar, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Philipp Koehn, and Christof Monz. Findings of the 2018 conference on machine translation (WMT18). In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 272–303, Belgium, Brussels, October 2018. Association for Computational Linguistics. doi: 10.18653/v1/W18-6401. URL <https://www.aclweb.org/anthology/W18-6401>.
- Raj Dabre, Chenhui Chu, and Anoop Kunchukuttan. A Brief Survey of Multilingual Neural Machine Translation. may 2019. URL <http://arxiv.org/abs/1905.05395>.
- Andreas Eisele and Yu Chen. {M}ulti{UN}: A Multilingual Corpus from United Nation Documents. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation ({LREC} ’10)*, Valletta, Malta, 2010. European Language Resources Association (ELRA). URL http://www.lrec-conf.org/proceedings/lrec2010/pdf/686{}_Paper.pdf.

- Lifeng Han. Machine Translation Evaluation Resources and Methods: A Survey. 2018, 2016. URL <http://arxiv.org/abs/1605.04515>.
- Melvin Johnson, Mike Schuster, Quoc V Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. Google’s Multilingual Neural Machine Translation System: Enabling Zero-Shot Translation. *Transactions of the Association for Computational Linguistics*, 5:339–351, 2017. doi: 10.1162/tacl.a_00065. URL <https://www.aclweb.org/anthology/Q17-1024>.
- Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F T Martins, and Alexandra Birch. Marian: Fast Neural Machine Translation in {C++}. In *Proceedings of ACL 2018, System Demonstrations*, Melbourne, Australia, 2018. URL <https://arxiv.org/abs/1804.00344>.
- Philipp Koehn. *Neural Machine Translation*. Cambridge University Press, 2020.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 311–318. Association for Computational Linguistics, 2002.
- Martin Popel and Ondřej Bojar. Training tips for the transformer model. *The Prague Bulletin of Mathematical Linguistics*, 110, 03 2018. doi: 10.2478/pralin-2018-0002.
- Aditya Siddhant, Melvin Johnson, Henry Tsai, Naveen Arivazhagan, Jason Riesa, Ankur Bapna, Orhan Firat, and Karthik Raman. Evaluating the Cross-Lingual Effectiveness of Massively Multilingual Neural Machine Translation. sep 2019. URL <http://arxiv.org/abs/1909.00437>.
- Jörg Tiedemann. Parallel data, tools and interfaces in opus. In Nicoletta Calzolari (Conference Chair), Khalid Choukri, Thierry Declerck, Mehmet Ugur Dogan, Bente Maegaard, Joseph Mariani, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC’12)*, Istanbul, Turkey, may 2012. European Language Resources Association (ELRA). ISBN 978-2-9517408-7-7.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In I Guyon, U V Luxburg, S Bengio, H Wallach, R Fergus, S Vishwanathan, and R Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf>.

List of Figures

1.1	Transformer model architecture	8
1.2	MNMT research categorized	9
1.3	Tranlsation performance for 102 languages from Arivazhagan et al. [2019]	9
2.1	Training data language statistics	12
2.2	Early stopping to prevent overfitting	15
2.3	Early stopping as the model is not improving	16
2.4	Single language results visualization: models other target languages are randomly selected vs. those selected from similar group of lan- guages	20
2.5	Overall model convergence dashboard	21
2.6	Example change of model's performance on validation set in time	22
2.7	Small improvement during long training	22
3.1	En→De bilingual results	23
3.2	En→De monolingual baseline results (RANDOM)	24
4.1	En→De BLEU score difference: Random vs. Germanic	28
4.2	En→Da BLEU score difference: Random vs. Germanic	29
4.3	En→Bg BLEU score difference: Random vs. Slavic with cyrillic script	30
4.4	En→Ru BLEU score difference: Random vs. Slavic with cyrillic script	31

List of Tables

1.1	Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types	4
1.2	BLEU scores for translation in one direction (part of Table 7 from [Aharoni et al., 2019])	7
2.1	Groups of subdatasets in OPUS	12
3.1	BLEU score change with adding target languages	25
A.1	BLEU scores for bilingual models	41

Abbreviations

ALPAC Automatic Language Processing Advisory Committee. 4

ARPA Advanced Research Projects Agency. 4

BLEU bilingual evaluation understudy. 4, 6, 9, 10, 13–15, 19, 20, 24, 26–30

CNN convolutional neural network. 3

GRU gated recurrent unit. 3

LSTM long short-term memory. 3

MT machine translation. 2, 5

NMT neural machine translation. 5, 6

RNN recurrent neural network. 3

SGD stochastic gradient descent. 13

Glossary

baseline In machine learning this term refers to a simple or naïve initial solution, which efficiency it then taken as a reference point and later improved. . 10

early stopping Regularization technique to avoid model overfit. Usually consists of stopping the training process when the value of some selected metric on the validation set is not improved for last number of validation steps. . 14, 15, 19

en-to-36 The dataset with source sentences in English and target sentences in 36 languages, described in Section 2.3.1 . 15, 16, 22

en-to-5 The dataset created from UN parallel corpus, with source sentences in English and target sentences in one of following 5 languages: Spanish, French, Russian, Arabic and Chinese; described in Section 2.3.2 . 15, 22

epoch Refers to one pass of full training dataset to the learning algorithm . 13

loss Loss function, also often called 'objective function' and 'error function'. It is optimized during the training process. In our experiments it is mean word cross-entropy score.. 19

overfitting Occurs when the model's performance on unseen validation set stops improving while on the training set it still improves. . 14, 20

A. Attachments

A.1 Additional tables

A.1.1 Bilingual results

target dataset	bg	da	de	is	mk	nl	no	ru	sv	uk
Books/v1	—	—	5.4	—	—	4.8	—	8.3	—	—
DGT/v4	33.1	27.4	24.1	—	—	25.9	—	—	28.9	—
ECB/v1	—	20.9	17.9	—	—	21.2	—	—	—	—
EMEA/v3	15.1	16.4	15.6	—	—	15.8	—	—	17.6	—
EUbookshop/v2	38.2	24.1	18.3	—	—	18.9	—	—	24.7	—
Europarl/v3	—	24.6	18.7	—	—	23.4	—	—	23.6	—
Europarl/v7	41.4	32.5	25.7	—	—	26.0	—	—	33.3	—
GNOME/v1	—	—	5.6	2.4	—	8.9	—	—	—	—
GlobalVoices/v2015	—	—	15.2	—	10.6	18.6	—	13.2	—	—
GlobalVoices/v2017q3	—	—	15.1	—	10.7	19.1	—	14.4	—	—
JRC-Acquis/v3.0	30.8	27.3	23.6	—	—	25.7	—	—	29.1	—
KDE4/v2	6.9	8.5	6.6	4.2	4.8	8.1	—	4.1	8.4	1.3
MPC1/v1	—	—	9.8	—	—	—	—	—	—	—
MultiUN/v1	—	—	25.4	—	—	—	—	14.6	—	—
News-Commentary/v11	—	—	18.4	—	—	19.2	—	23.9	—	—
News-Commentary/v9.0	—	—	13.2	—	—	—	—	18.2	—	—
News-Commentary/v9.1	—	—	19.3	—	—	—	—	22.1	—	—
OpenOffice/v2	—	—	8.7	—	—	—	—	—	8.6	—
OpenSubtitles/v1	19.3	17.1	10.8	—	—	12.5	23.1	16.2	13.4	—
OpenSubtitles/v2016	23.2	14.8	13.0	24.3	24.3	13.7	27.0	19.5	14.8	11.2
OpenSubtitles/v2018	23.7	15.6	13.1	23.1	24.6	13.4	29.6	19.2	15.3	12.2
PHP/v1	—	—	7.2	—	—	12.6	—	3.3	8.9	—
ParaCrawl/v1	—	—	12.5	—	—	17.9	—	11.1	—	—
SETIMES/v1	23.2	—	—	—	6.4	—	—	—	—	—
SETIMES/v2	27.5	—	—	—	10.4	—	—	—	—	—
TED2013/v1.1	—	—	16.9	—	—	19.1	—	14.7	—	—
Tanzil/v1	5.7	—	5.8	—	—	4.6	6.1	2.4	4.4	—
Tatoeba/v2	—	—	22.6	—	—	28.9	—	27.7	—	13.3
UN/v20090831	—	—	—	—	—	—	—	9.9	—	—
Ubuntu/v14.10	—	—	—	—	—	8.6	—	—	—	—
WMT-News/v2014	—	—	13.9	—	—	—	—	—	—	—
WikiSource/v1	—	—	—	—	—	—	—	—	5.3	—
Wikipedia/v1.0	11.7	—	7.8	—	—	9.6	—	10.6	—	—

Table A.1: BLEU scores for bilingual models