

# Fully Character-Level Neural Machine Translation without Explicit Segmentation

**Jason Lee\***  
ETH Zürich  
jasonlee@inf.ethz.ch

**Kyunghyun Cho**  
New York University  
kyunghyun.cho@nyu.edu

**Thomas Hofmann**  
ETH Zürich  
thomas.hofmann@inf.ethz.ch

## Abstract

Most existing machine translation systems operate at the level of words, relying on explicit segmentation to extract tokens. We introduce a neural machine translation (NMT) model that maps a source character sequence to a target character sequence without any segmentation. We employ a character-level convolutional network with max-pooling at the encoder to reduce the length of source representation, allowing the model to be trained at a speed comparable to subword-level models while capturing local regularities. Our character-to-character model outperforms a recently proposed baseline with a subword-level encoder on WMT'15 DE-EN and CS-EN, and gives comparable performance on FI-EN and RU-EN. We then demonstrate that it is possible to share a single character-level encoder across multiple languages by training a model on a many-to-one translation task. In this multilingual setting, the character-level encoder significantly outperforms the subword-level encoder on all the language pairs. We observe that on CS-EN, FI-EN and RU-EN, the quality of the multilingual character-level translation even surpasses the models specifically trained on that language pair alone, both in terms of the BLEU score and human judgment.

## 1 Introduction

Nearly all previous work in machine translation has been at the level of words. Aside from our intu-

\*The majority of this work was completed while the author was visiting New York University.

itive understanding of word as a basic unit of meaning (Jackendoff, 1992), one reason behind this is that sequences are significantly longer when represented in characters, compounding the problem of data sparsity and modeling long-range dependencies. This has driven NMT research to be almost exclusively word-level (Bahdanau et al., 2015; Sutskever et al., 2014).

Despite their remarkable success, word-level NMT models suffer from several major weaknesses. For one, they are unable to model rare, out-of-vocabulary words, making them limited in translating languages with rich morphology such as Czech, Finnish and Turkish. If one uses a large vocabulary to combat this (Jean et al., 2015), the complexity of training and decoding grows linearly with respect to the target vocabulary size, leading to a vicious cycle.

To address this, we present a fully character-level NMT model that maps a character sequence in a source language to a character sequence in a target language. We show that our model outperforms a baseline with a subword-level encoder on DE-EN and CS-EN, and achieves a comparable result on FI-EN and RU-EN. A purely character-level NMT model with a basic encoder was proposed as a baseline by Luong and Manning (2016), but training it was prohibitively slow. We were able to train our model at a reasonable speed by drastically reducing the length of source sentence representation using a stack of convolutional, pooling and highway layers.

One advantage of character-level models is that they are better suited for multilingual translation than their word-level counterparts which require a separate word vocabulary for each language. We

verify this by training a single model to translate four languages (German, Czech, Finnish and Russian) to English. Our multilingual character-level model outperforms the subword-level baseline by a considerable margin in all four language pairs, strongly indicating that a character-level model is more flexible in assigning its capacity to different language pairs. Furthermore, we observe that our multilingual character-level translation even exceeds the quality of bilingual translation in three out of four language pairs, both in BLEU score metric and human evaluation. This demonstrates excellent parameter efficiency of character-level translation in a multilingual setting. We also showcase our model’s ability to handle intra-sentence code-switching while performing language identification on the fly.

The contributions of this work are twofold: we empirically show that (1) we can train character-to-character NMT model without any explicit segmentation; and (2) we can share a single character-level encoder across multiple languages to build a multilingual translation system without increasing the model size.

## 2 Background: Attentional Neural Machine Translation

Neural machine translation (NMT) is a recently proposed approach to machine translation that builds a single neural network which takes as an input, a source sentence  $X = (x_1, \dots, x_{T_X})$  and generates its translation  $Y = (y_1, \dots, y_{T_Y})$ , where  $x_t$  and  $y_{t'}$  are source and target symbols (Bahdanau et al., 2015; Sutskever et al., 2014; Luong et al., 2015; Cho et al., 2014a). Attentional NMT models have three components: an *encoder*, a *decoder* and an *attention* mechanism.

**Encoder** Given a source sentence  $X$ , the encoder constructs a continuous representation that summarizes its meaning with a recurrent neural network (RNN). A bidirectional RNN is often implemented as proposed in (Bahdanau et al., 2015). A forward encoder reads the input sentence from left to right:  $\vec{\mathbf{h}}_t = \vec{f}_{\text{enc}}(E_x(x_t), \vec{\mathbf{h}}_{t-1})$ . Similarly, a backward encoder reads it from right to left:  $\overleftarrow{\mathbf{h}}_t = \overleftarrow{f}_{\text{enc}}(E_x(x_t), \overleftarrow{\mathbf{h}}_{t+1})$ , where  $E_x$  is

the source embedding lookup table, and  $\vec{f}_{\text{enc}}$  and  $\overleftarrow{f}_{\text{enc}}$  are recurrent activation functions such as long short-term memory units (LSTMs) (Hochreiter and Schmidhuber, 1997) or gated recurrent units (GRUs) (Cho et al., 2014b). The encoder constructs a set of continuous source sentence representations  $C$  by concatenating the forward and backward hidden states at each timestep:  $C = \{\mathbf{h}_1, \dots, \mathbf{h}_{T_X}\}$ , where  $\mathbf{h}_t = [\vec{\mathbf{h}}_t; \overleftarrow{\mathbf{h}}_t]$ .

**Attention** First introduced in Bahdanau et al. (2015), the attention mechanism lets the decoder *attend* more to different source symbols for each target symbol. More concretely, it computes the context vector  $\mathbf{c}_{t'}$  at each decoding time step  $t'$  as a weighted sum of the source hidden states:  $\mathbf{c}_{t'} = \sum_{t=1}^{T_X} \alpha_{t't} \mathbf{h}_t$ . Similarly to Chung et al. (2016) and Firat et al. (2016a), each attentional weight  $\alpha_{t't}$  represents how relevant the  $t$ -th source token  $x_t$  is to the  $t'$ -th target token  $y_{t'}$ , and is computed as:

$$\alpha_{t't} = \frac{1}{Z} \exp\left(\text{score}\left(E_y(y_{t'-1}), \mathbf{s}_{t'-1}, \mathbf{h}_t\right)\right), \quad (1)$$

where  $Z = \sum_{k=1}^{T_X} \exp(\text{score}(E_y(y_{t'-1}), \mathbf{s}_{t'-1}, \mathbf{h}_k))$  is the normalization constant.  $\text{score}()$  is a feed-forward neural network with a single hidden layer that scores how well the source symbol  $x_t$  and the target symbol  $y_{t'}$  match.  $E_y$  is the target embedding lookup table and  $\mathbf{s}_{t'}$  is the target hidden state at time  $t'$ .

**Decoder** Given a source context vector  $\mathbf{c}_{t'}$ , the decoder computes its hidden state at time  $t'$  as:  $\mathbf{s}_{t'} = f_{\text{dec}}(E_y(y_{t'-1}), \mathbf{s}_{t'-1}, \mathbf{c}_{t'})$ . Then, a parametric function  $\text{out}_k()$  returns the conditional probability of the next target symbol being  $k$ :

$$p(y_{t'} = k | y_{<t'}, X) = \frac{1}{Z} \exp\left(\text{out}_k\left(E_y(y_{t'-1}), \mathbf{s}_{t'}, \mathbf{c}_{t'}\right)\right) \quad (2)$$

where  $Z$  is again the normalization constant:  $Z = \sum_j \exp(\text{out}_j(E_y(y_{t'-1}), \mathbf{s}_{t'}, \mathbf{c}_{t'}))$ .

**Training** The entire model can be trained end-to-end by minimizing the negative conditional log-

likelihood, which is defined as:

$$\mathcal{L} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_Y^{(n)}} \log p(y_t = y_t^{(n)} | y_{<t}^{(n)}, X^{(n)}),$$

where  $N$  is the number of sentence pairs, and  $X^{(n)}$  and  $y_t^{(n)}$  are the source sentence and the  $t$ -th target symbol in the  $n$ -th pair, respectively.

### 3 Fully Character-Level Translation

#### 3.1 Why Character-Level?

The benefits of character-level translation over word-level translation are well known. Chung et al. (2016) present three main arguments: character level models (1) do not suffer from out-of-vocabulary issues, (2) are able to model different, rare morphological variants of a word, and (3) do not require segmentation. Particularly, text segmentation is highly non-trivial for many languages and problematic even for English as word tokenizers are either manually designed or trained on a corpus using an objective function that is unrelated to the translation task at hand, which makes the overall system sub-optimal.

Here we present two additional arguments for character-level translation. First, a character-level translation system can easily be applied to a multilingual translation setting. Between European languages where the majority of alphabets overlaps, for instance, a character-level model may easily identify morphemes that are shared across different languages. A word-level model, however, will need a separate word vocabulary for each language, allowing no cross-lingual parameter sharing.

Also, by not segmenting source sentences into words, we no longer inject our knowledge of words and word boundaries into the system; instead, we encourage the model to discover an internal structure of a sentence by itself and learn how a sequence of symbols can be mapped to a continuous meaning representation.

#### 3.2 Related Work

To address these limitations associated with word-level translation, a recent line of research has investigated using sub-word information.

Costa-Jussá and Fonollosa (2016) replaced the word-lookup table with convolutional and highway

layers on top of character embeddings, while still segmenting source sentences into words. Target sentences were also segmented into words, and predictions were made at word-level.

Similarly, Ling et al. (2015) employed a bidirectional LSTM to compose character embeddings into word embeddings. At the target side, another LSTM takes the hidden state of the decoder and generates the target word, character by character. While this system is completely open-vocabulary, it also requires offline segmentation. Character-to-word and word-to-character LSTMs significantly slow down training, as well.

Most recently, Luong and Manning (2016) proposed a hybrid scheme that consults character-level information whenever the model encounters an out-of-vocabulary word. As a baseline, they also implemented a purely character-level NMT model with 4 layers of unidirectional LSTMs with 512 cells, with attention over each character. Despite being extremely slow (approximately 3 months to train), the character-level model gave a comparable performance to the word-level baseline. This shows the possibility of fully character-level translation.

Having a word-level decoder restricts the model to only being able to generate previously seen words. Sennrich et al. (2015) introduced a subword-level NMT model that is capable of open-vocabulary translation using subword-level segmentation based on the byte pair encoding (BPE) algorithm. Starting from a character vocabulary, the algorithm identifies frequent character  $n$ -grams in the training data and iteratively adds them to the vocabulary, ultimately giving a subword vocabulary which consists of words, subwords and characters. Once the segmentation rules have been learned, their model performs subword-to-subword translation (**bpe2bpe**) in the same way as word-to-word translation.

Perhaps the work that is closest to our end goal is (Chung et al., 2016), which used a subword-level encoder from (Sennrich et al., 2015) and a fully character-level decoder (**bpe2char**). Their results show that character-level decoding performs better than subword-level decoding. Motivated by this work, we aim for fully character-level translation at both sides (**char2char**).

Outside NMT, our work is based on a few existing approaches that applied convolutional networks

to text, most notably in text classification (Zhang et al., 2015; Xiao and Cho, 2016). We also drew inspiration for our multilingual models from previous work that showed the possibility of training a single recurrent model for multiple languages in domains other than translation (Tsvetkov et al., 2016; Gillick et al., 2015).

### 3.3 Challenges

Sentences are on average 6 (DE, CS and RU) to 8 (FI) times longer when represented in characters. This poses three major challenges to achieving fully character-level translation.

**(1) Training/decoding latency** For the decoder, although the sequence to be generated is much longer, each character-level softmax operation costs considerably less compared to a word- or subword-level softmax. Chung et al. (2016) report that character-level decoding is only 14% slower than subword-level decoding.

On the other hand, computational complexity of the attention mechanism grows quadratically with respect to the sentence length, as it needs to attend to every source token for every target token. This makes a naive character-level approach, such as in Luong and Manning (2016), computationally prohibitive. Consequently, reducing the length of the source sequence is key to ensuring reasonable speed in both training and decoding.

**(2) Mapping character sequence to continuous representation** The arbitrary relationship between the orthography of a word and its meaning is a well-known problem in linguistics (de Saussure, 1916). Building a character-level encoder is arguably a more difficult problem, as the encoder needs to learn a highly non-linear function from a long sequence of character symbols to a meaning representation.

**(3) Long range dependencies in characters** A character-level encoder needs to model dependencies over longer timespans than a word-level encoder does.

## 4 Fully Character-Level NMT

### 4.1 Encoder

We design an encoder that addresses all the challenges discussed above by using convolutional and pooling layers aggressively to both (1) drastically shorten the input sentence; and (2) efficiently capture local regularities. Inspired by the character-level language model from Kim et al. (2015), our encoder first reduces the source sentence length with a series of convolutional, pooling and highway layers. The shorter representation, instead of the full character sequence, is passed through a bidirectional GRU to (3) help it resolve long term dependencies. We illustrate the proposed encoder in Figure 1 and discuss each layer in detail below.

**Embedding** We map the sequence of source characters  $(x_1, \dots, x_{T_x})$  to a sequence of character embeddings of dimensionality  $d_c$ :  $X = (\mathbf{C}(x_1), \dots, \mathbf{C}(x_{T_x})) \in \mathbb{R}^{d_c \times T_x}$  where  $T_x$  is the number of source characters and  $\mathbf{C}$  is the character embedding lookup table:  $\mathbf{C} \in \mathbb{R}^{d_c \times |\mathcal{C}|}$ .

**Convolution** One-dimensional convolution operation is then used along consecutive character embeddings. Assuming we have a single filter  $\mathbf{f} \in \mathbb{R}^{d_c \times w}$  of width  $w$ , we first apply padding to the beginning and the end of  $X$ , such that the padded sentence  $X' \in \mathbb{R}^{d_c \times (T_x + w - 1)}$  is  $w - 1$  symbols longer. We then apply a narrow convolution between  $X'$  and  $\mathbf{f}$  such that the  $k$ -th element of the output  $Y_k$  is given as:

$$Y_k = (X' * \mathbf{f})_k = \sum_{i,j} (X'_{[:,k-w+1:k]} \otimes \mathbf{f})_{ij}, \quad (3)$$

where  $\otimes$  denotes elementwise matrix multiplication and  $*$  is the convolution operation.  $X'_{[:,k-w+1:k]}$  is the sliced subset of  $X'$  that contains all the rows but only  $w$  adjacent columns. The padding scheme employed above, commonly known as *half convolution*, ensures that the length of the output is identical to the length of the input, (i.e.,  $Y \in \mathbb{R}^{1 \times T_x}$ ).

We just illustrated how a single convolutional filter of fixed width might be applied to a sentence. In order to extract informative character patterns of different lengths, we employ a set of filters of varying widths. More concretely, we use a filter

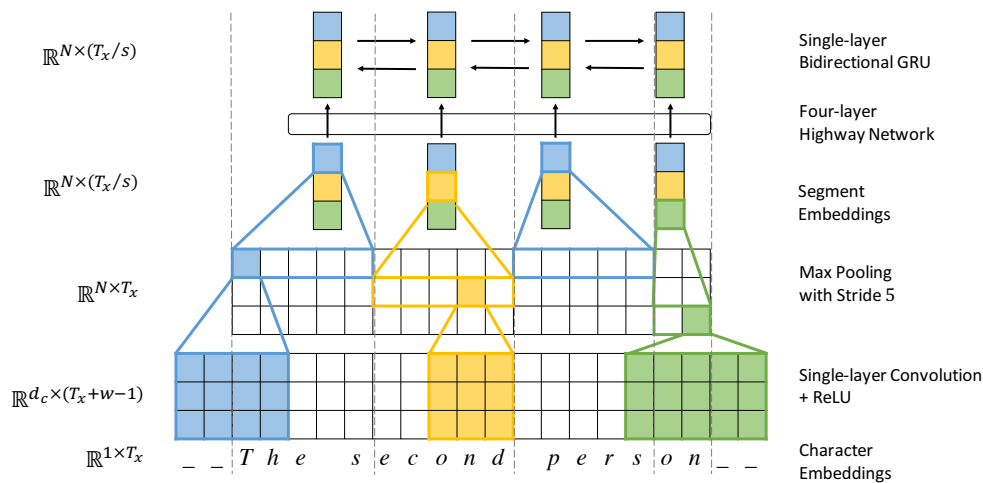


Figure 1: Encoder architecture schematics. Underscore denotes padding. A dotted vertical line delimits each segment. The stride of pooling  $s$  is 5 in the diagram.

bank  $\mathbf{F} = \{\mathbf{f}_1, \dots, \mathbf{f}_m\}$  where  $\mathbf{f}_i = \mathbb{R}^{d_c \times i \times n_i}$  is a collection of  $n_i$  filters of width  $i$ . Our model uses  $m = 8$ , hence extracting character n-grams up to 8 characters long. Outputs from all the filters are stacked upon each other, giving a single representation  $Y \in \mathbb{R}^{N \times T_x}$ , where the dimensionality of each column is given by the total number of filters  $N = \sum_{i=1}^m n_i$ . Finally, rectified linear activation (ReLU) is applied elementwise to this representation.

**Max pooling with stride** The output from the convolutional layer is first split into segments of width  $s$ , and max-pooling over time is applied to each segment with no overlap. This procedure selects the most salient features to give a *segment embedding*. Each segment embedding is a summary of meaningful character n-grams occurring in a particular (overlapping) subsequence in the source sentence. Note that the rightmost segment (above ‘on’) in Figure 1 may capture ‘son’ (the filter in green) although ‘s’ occurs in the previous segment. In other words, our segments are overlapping as opposed to in word- or subword-level models with hard segmentation.

Segments act as our internal linguistic unit from this layer and above: the attention mechanism, for instance, attends to each source segment instead of source character. This shortens the source representation  $s$ -fold:  $Y' \in \mathbb{R}^{N \times (T_x/s)}$ . Empirically, we found using a smaller  $s$  leads to better performance

at increased training time. We chose  $s = 5$  in our experiments as it gives a reasonable balance between the two.

**Highway network** A sequence of segment embeddings from the max pooling layer is fed into a highway network (Srivastava et al., 2015). Highway networks are shown to significantly improve the quality of a character-level language model when used with convolutional layers (Kim et al., 2015). A highway network transforms input  $\mathbf{x}$  with a gating mechanism that adaptively regulates information flow:

$$\mathbf{y} = g \odot \text{ReLU}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + (1 - g) \odot \mathbf{x},$$

where  $g = \sigma((\mathbf{W}_2 \mathbf{x} + \mathbf{b}_2))$ . We apply this to each segment embedding individually.

**Recurrent layer** Finally, the output from the highway layer is given to a bidirectional GRU from §2, using each segment embedding as input.

**Subword-level encoder** Unlike a subword-level encoder, our model does not commit to a specific choice of segmentation; instead it is trained to consider every possible character pattern and extract only the most meaningful ones. Therefore, the definition of segmentation in our model is dynamic unlike subword-level encoders. During training, the model finds the most salient character patterns in a sentence via max-pooling, and the character

<b>Bilingual</b>	bpe2char	char2char
Vocab size	24,440	300
Source emb.	512	128
Target emb.	512	512
Conv. filters		200-200-250-250 -300-300-300-300
Pool stride		5
Highway		4 layers
Encoder	1-layer 512 GRUs	
Decoder	2-layer 1024 GRUs	

Table 1: Bilingual model architectures. The char2char model uses 200 filters of width 1, 200 filters of width 2, ... and 300 filters of width 8.

sequences extracted by the model change over the course of training. This is in contrast to how BPE segmentation rules are learned: the segmentation is learned and fixed before training begins.

## 4.2 Attention and Decoder

Similarly to the attention model in Chung et al. (2016) and Firat et al. (2016a), a single-layer feed-forward network computes the attention score of next target character to be generated with every source segment representation. A standard two-layer character-level decoder then takes the source context vector from the attention mechanism and predicts each target character. This decoder was described as *base decoder* by Chung et al. (2016).

## 5 Experiment Settings

### 5.1 Task and Models

We evaluate the proposed character-to-character (**char2char**) translation model against subword-level baselines (**bpe2bpe** and **bpe2char**) on the WMT’15 DE→EN, CS→EN, FI→EN and RU→EN translation tasks.<sup>1</sup> We do not consider word-level models, as it has already been shown that subword-level models outperform them by mitigating issues inherent to closed-vocabulary translation (Sennrich et al., 2015; Sennrich et al., 2016). Indeed, subword-level NMT models have been the de-facto state-of-the-art and are now used in a very large-scale industry NMT system to serve millions of users per day (Wu et al., 2016).

<sup>1</sup><http://www.statmt.org/wmt15/translation-task.html>

We experiment in two different scenarios: 1) a bilingual setting where we train a model on data from a single language pair; and 2) a multilingual setting where the task is many-to-one translation. We train a single model on data from all four language pairs. Hence, our baselines and models are:

- (a) bilingual bpe2bpe: from (Firat et al., 2016a)
- (b) bilingual bpe2char: from (Chung et al., 2016)
- (c) bilingual char2char
- (d) multilingual bpe2char
- (e) multilingual char2char

We train all the models ourselves other than (a), for which we report the results from Firat et al. (2016a). We detail the configuration of our models in Table 1 and Table 2.

### 5.2 Datasets and Preprocessing

We use all available parallel data on the four language pairs from WMT’15: DE-EN, CS-EN, FI-EN and RU-EN.

For the bpe2char baselines, we only use sentence pairs where the source is no longer than 50 subword symbols. For our char2char models, we only use pairs where the source sentence is no longer than 450 characters. For all the language pairs apart from FI-EN, we use newstest-2013 as a development set and newstest-2014 and newstest-2015 as test sets. For FI-EN, we use newsdev-2015 and newstest-2015 as development and test sets, respectively. We tokenize<sup>2</sup> each corpus using the script from Moses.<sup>3</sup>

When training bilingual bpe2char models, we extract 20,000 BPE operations from each of the source and target corpus using a script from Sennrich et al. (2015). This gives a source BPE vocabulary of size 20k–24k for each language.

### 5.3 Training Details

Each model is trained using stochastic gradient descent and Adam (Kingma and Ba, 2014) with a learning rate of 0.0001 and minibatch size 64. Training continues until the BLEU score on the validation

<sup>2</sup>This is unnecessary for char2char models, yet was carried out for comparison.

<sup>3</sup><https://github.com/moses-smt/mosesdecoder>

Multilingual	bpe2char	char2char
Vocab size	54,544	400
Source emb.	512	128
Target emb.	512	512
Conv. filters		200-250-300-300 -400-400-400-400
Pool stride		5
Highway		4 layers
Encoder	1-layer 512 GRUs	
Decoder	2-layer 1024 GRUs	

Table 2: Multilingual model architectures.

set stops improving. The norm of the gradient is clipped with a threshold of 1 (Pascanu et al., 2013). All weights are initialized from a uniform distribution  $[-0.01, 0.01]$ .

Each model is trained on a single pre-2016 GTX Titan X GPU with 12GB RAM.

#### 5.4 Decoding Details

As done by Chung et al. (2016), a two-layer unidirectional character-level decoder with 1024 GRU units is used for all our experiments. For decoding, we use a beam search algorithm with length-normalization to penalize shorter hypotheses. The beam width is 20 for all models.

#### 5.5 Training Multilingual Models

**Task description** We train a model on a many-to-one translation task to translate a sentence in any of the four languages (German, Czech, Finnish and Russian) to English. We do *not* provide a language identifier to the encoder, but merely the sentence itself, encouraging the model to perform language identification on the fly. In addition, by not providing the language identifier, we expect the model to handle intra-sentence code-switching seamlessly.

**Model architecture** The multilingual char2char model uses slightly more convolutional filters than the bilingual char2char model, namely (200-250-300-300-400-400-400-400). Otherwise, the architecture remains the same as shown in Table 1. By not changing the size of the encoder and the decoder, we fix the capacity of the core translation module, and only allow the multilingual model to detect more character patterns.

Similarly, the multilingual bpe2char model has the same encoder and decoder as the bilingual bpe2char model, but a larger vocabulary. We learn 50,000 multilingual BPE operations on the multilingual corpus, resulting in 54,544 subwords. See Table 2 for the exact configuration of our multilingual models.

**Data scheduling** For the multilingual models, an appropriate scheduling of data from different languages is crucial to avoid overfitting to one language too soon. Following Firat et al. (2016a) and Firat et al. (2016b), each minibatch is *balanced*, in that the proportion of each language pair in a single minibatch corresponds to that of the full corpus. With this minibatch scheme, roughly the same number of updates is required to make one full pass over the entire training corpus of each language pair. Minibatches from all language pairs are combined and presented to the model as a single minibatch. See Table 3 for the minibatch size for each language pair.

	DE-EN	CS-EN	FI-EN	RU-EN
corpus size	4.5m	12.1m	1.9m	2.3m
minibatch size	14	37	6	7

Table 3: The minibatch size of each language (second row) is proportionate to the number of sentence pairs in each corpus (first row).

**Treatment of Cyrillic** To facilitate cross-lingual parameter sharing, we convert every Cyrillic character in the Russian source corpus to Latin alphabet according to ISO-9. Table 4 shows an example of how this conversion may help the multilingual models identify lexemes that are shared across multiple languages.

	school	schools
CS	škola	školy
RU	школа	школы
RU (ISO-9)	škola	školy

Table 4: Czech and Russian words for *school* and *schools*, alongside the conversion of Russian characters into Latin.

**Multilingual BPE** For the multilingual bpe2char model, multilingual BPE segmentation rules are extracted from a large dataset containing training source corpora of all the language pairs. To ensure the BPE rules are not biased towards one language,

	Setting	Src	Trg	Dev	Test1	Test2	
DE-EN	(a)*	bi	bpe	bpe	24.13		24.00
	(b)	bi	bpe	char	25.64	24.59	25.27
	(c)	bi	char	char	<b>26.30</b>	<b>25.77</b>	<b>25.83</b>
	(d)	multi	bpe	char	24.92	24.54	25.23
	(e)	multi	char	char	25.67	25.13	25.79
CS-EN	(f)*	bi	bpe	bpe	21.24		20.32
	(g)	bi	bpe	char	22.95	23.78	22.40
	(h)	bi	char	char	23.38	24.08	22.46
	(i)	multi	bpe	char	23.27	24.27	22.42
	(j)	multi	char	char	<b>24.09</b>	<b>25.01</b>	<b>23.24</b>
FI-EN	(k)*	bi	bpe	bpe	13.15		12.24
	(l)	bi	bpe	char	14.54		13.98
	(m)	bi	char	char	14.18		13.10
	(n)	multi	bpe	char	14.70		14.40
	(o)	multi	char	char	<b>15.96</b>		<b>15.74</b>
RU-EN	(p)*	bi	bpe	bpe	21.04		22.44
	(q)	bi	bpe	char	21.68	26.21	22.83
	(r)	bi	char	char	21.75	<b>26.80</b>	22.73
	(s)	multi	bpe	char	21.75	26.31	22.81
	(t)	multi	char	char	<b>22.20</b>	26.33	<b>23.33</b>

Table 5: BLEU scores of five different models on four language pairs. For each test or development set, the best performing model is shown in bold. (\*) results are taken from (Firat et al., 2016a).

larger datasets such as Czech and German corpora are trimmed such that every corpus contains, approximately, an equal number of characters.

## 6 Quantitative Analysis

### 6.1 Evaluation with BLEU Score

In this section, we first establish our main hypotheses for introducing character-level and multilingual models, and investigate whether our observations support or disagree with our hypotheses. From our empirical results, we want to verify: (1) if fully character-level translation outperforms subword-level translation, (2) in which setting and to what extent is multilingual translation beneficial and (3) if multilingual, character-level translation achieves superior performance to other models. We outline our results with respect to each hypothesis below.

**(1) Character-level vs. subword-level** In a bilingual setting, the char2char model outperforms both subword-level baselines on DE-EN (Table 5 (a-c)) and CS-EN (Table 5 (f-h)). On the other two language pairs, it exceeds the bpe2bpe model and achieves a similar performance with the bpe2char baseline (Table 5 (k-m) and (p-r)). We conclude that

the proposed character-level model is comparable or better than both subword-level baselines.

Meanwhile, in a multilingual setting, the character-level encoder significantly surpasses the subword-level encoder consistently in all the language pairs (Table 5 (d-e), (i-j), (n-o) and (s-t)). From this, we conclude that translating at the level of characters allows the model to discover shared constructs between languages more effectively. This also demonstrates that the character-level model is more flexible in assigning model capacity to different language pairs.

**(2) Multilingual vs. bilingual** At the level of characters, we note that multilingual translation is indeed strongly beneficial. On the test sets, the multilingual character-level model outperforms the single-pair character-level model by 2.64 BLEU in FI-EN (Table 5 (m, o)) and 0.78 BLEU in CS-EN (Table 5 (h, j)), while achieving comparable results on DE-EN and RU-EN.

At the level of subwords, on the other hand, we do not observe the same degree of performance benefit. The multilingual bpe2char model requires much more updates to reach the performance of the bilingual bpe2char model (see Figure 2). This



	Setting	Src	Trg	Adequacy		Fluency		
				Raw (%)	Std. ( $\sigma$ )	Raw (%)	Std. ( $\sigma$ )	
DE-EN	(a)	bi	bpe	char	65.47	-0.0536	68.64	0.0052
	(b)	bi	char	char	68.11	<b>0.0509</b>	68.80	0.0468
	(c)	multi	char	char	67.80	<b>0.0281</b>	68.92	0.0282
CS-EN	(d)	bi	bpe	char	62.76	<b>0.0361</b>	61.62	-0.0285
	(e)	bi	char	char	60.78	-0.0154	63.37	0.0410
	(f)	multi	char	char	63.03	<b>0.0415</b>	65.08	<b>0.1047</b>
FI-EN	(g)	bi	bpe	char	47.03	-0.1326	59.33	-0.0329
	(h)	bi	char	char	50.17	<b>-0.0650</b>	59.97	-0.0216
	(i)	multi	char	char	50.95	<b>-0.0110</b>	63.26	<b>0.0969</b>
RU-EN	(j)	bi	bpe	char	61.26	-0.1062	57.74	-0.0592
	(k)	bi	char	char	64.06	<b>0.0105</b>	59.85	0.0168
	(l)	multi	char	char	64.77	<b>0.0116</b>	63.32	<b>0.1748</b>

Table 6: Human evaluation results for adequacy and fluency. We present both the averaged raw scores (Raw) and the averaged standardized scores (Std.). Standardized adequacy is used to rank the systems and standardized fluency is used to break ties. A positive standardized score should be interpreted as the number of standard deviations above this particular worker’s mean score that this system scored on average. For each language pair, we boldface the best performing model with statistical significance. When there is a tie, we boldface both systems.

suggests that learning useful subword segmentation across languages is difficult.

**(3) Multilingual char2char vs. others** The multilingual char2char model is the best performer in CS-EN, FI-EN and RU-EN (Table 5 (j, o, t)), and is the runner-up in DE-EN (Table 5 (e)). The fact that the multilingual char2char model outperforms the single-pair models goes to show the parameter efficiency of character-level translation: instead of training  $N$  separate models for  $N$  language pairs, it is possible to get a better performance with a single multilingual character-level model.

## 6.2 Human Evaluation

It is well known that automatic evaluation metrics such as BLEU encourage reference-like translations and do not fully capture true translation quality (Callison-Burch, 2009; Graham et al., 2015). Therefore, we also carry out a recently proposed evaluation from Graham et al. (2017) where we have human assessors rate both (1) adequacy; and (2) fluency of each system translation on a scale from 0 to 100 via Amazon Mechanical Turk. Adequacy is the degree to which assessors agree that the system translation expresses the meaning of the reference translation. Fluency is evaluated using system translation alone without any reference translation.

Approximately 1K Turkers assessed a single test set (3K sentences in newstest-2014) for each system and language pair. Each Turker conducted a minimum of 100 assessments for quality control, and the set of scores generated by each Turker was standardized to remove any bias in the individual’s scoring strategy.

We consider three models (bilingual bpe2char, bilingual char2char and multilingual char2char) for the human evaluation. We leave out the multilingual bpe2char model to minimize the number of similar systems to improve the interpretability of the evaluation overall.

For DE-EN, we observe that the multilingual char2char and bilingual char2char models are tied with respect to both adequacy and fluency (Table 6 (b-c)). For CS-EN, the multilingual char2char and bilingual bpe2char models are tied for adequacy. However, the multilingual char2char model yields significantly better fluency (Table 6 (d, f)). For FI-EN and RU-EN, the multilingual char2char model is tied with the bilingual char2char model with respect to adequacy, but significantly outperforms all other models in fluency (Table 6 (g-i, j-l)).

Overall, the improvement in translation quality yielded by the multilingual character-level model mainly comes from fluency. We conjecture that because the English decoder of the multilingual model is tuned in on all the training sentence pairs, it

**(a) Spelling mistakes**

DE ori	Warum sollten wir nicht Freunde sei ?
DE src	Warum solltne wir nich Freunde sei ?
EN ref	Why should not we be friends ?
bpe2char	Why are we to be friends ?
char2char	Why should we not be friends ?

**(b) Rare words**

DE src	Siebentausendzweihundertvierundfünfzig .
EN ref	Seven thousand two hundred fifty four .
bpe2char	Fifty-five Decline of the Seventy .
char2char	Seven thousand hundred thousand fifties .

**(c) Morphology**

DE src	Die Zufahrtsstraßen wurden gesperrt , wodurch sich laut CNN lange Rückstaus bildeten .
EN ref	The access roads were blocked off , which , according to CNN , caused long tailbacks .
bpe2char	The access roads were locked , which , according to CNN , was long back .
char2char	The access roads were blocked , which looked long backwards , according to CNN .

**(d) Nonce words**

DE src	Der Test ist nun über , aber ich habe keine gute Note . Es ist wie eine Verschlimmbesserung .
EN ref	The test is now over , but i don't have any good grade . it is like a worsened improvement .
bpe2char	The test is now over , but i do not have a good note .
char2char	The test is now , but i have no good note , it is like a worsening improvement .

**(e) Multilingual**

Multi src	Bei der <b>Metropolitního výboru pro dopravu</b> für das Gebiet der San Francisco Bay erklärten Beamte , der <b>Kongress könne das Problem</b> <b>банкротство доверительного Фонда строительства шоссеиных дорог</b> einfach durch Erhöhung der Kraftstoffsteuer lösen .
EN ref	At the Metropolitan Transportation Commission in the San Francisco Bay Area , officials say Congress could very simply <b>deal with the bankrupt Highway Trust Fund</b> by raising gas taxes .
bpe2char	During the Metropolitan Committee on Transport for San Francisco Bay , officials declared that Congress could <b>solve the problem of bankruptcy</b> by increasing the fuel tax bankrupt .
char2char	At the Metropolitan Committee on Transport for the territory of San Francisco Bay , officials explained that the Congress could simply <b>solve the problem of the bankruptcy of the Road Construction Fund</b> by increasing the fuel tax .

Table 7: Sample translations. For each example, we show the source sentence as *src*, the human translation as *ref*, and the translations from the subword-level baseline and our character-level model as *bpe2char* and *char2char*, respectively. For (a), the original, uncorrupted source sentence is also shown (*ori*). The source sentence in (e) contains words in German (in green), Czech (in yellow) and Russian (in blue). The translations in (a-d) are from the bilingual models, whereas those in (e) are from the multilingual models.

becomes a better language model than a bilingual model’s decoder. We leave it for future work to confirm if this is indeed the case.

## 7 Qualitative Analysis

In Table 7, we demonstrate our character-level model’s robustness in four translation scenarios from which conventional NMT systems are known to suffer. We also showcase our model’s ability to seamlessly handle intra-sentence *code-switching*, or mixed utterances from two or more languages.

We compare sample translations from the character-level model with those from the subword-level model, which already sidesteps some of the issues associated with word-level translation.

With real-world text containing typos and spelling mistakes, the quality of word-based translation would severely drop, as every non-canonical form of a word cannot be represented. On the other hand, a character-level model has a much better chance recovering the original word or sentence. Indeed, our char2char model is robust against a few spelling

(f) Long-distance dependencies

DE src	Der Rückgang zusammen mit einem verstärkten Sinken der Anzahl der Hausbesitzer unter 35 Jahren könnte dazu führen , dass Gartenzentren zehntausende Pfund pro Jahr verlieren , wenn die heutigen <b>jungen Konsumenten</b> nach einer Studie der HTA , wie in der Financial Times berichtet , die “ <b>Kernaltersgruppe für Gartenprodukte</b> ” erreichen .
EN ref	The drop , coupled with a particular decline in the number of homeowners aged under 35 , could result in garden centres losing out on tens of millions of pounds a year when today ’s <b>young consumers reach the “ core gardening age group , ”</b> according to the HTA ’s study , which was reported by the Financial Times .
bpe2char	The decline , together with reinforcing sinks of the number of householders under the age of 35 , could lead to tens of thousands of Garden Centres losing tens of thousands of pounds a year if today ’s <b>young consumers reach the “ kernel group of gardening products ”</b> according to a study of the HTA , as reported in the Financial Times .
char2char	The decline , together with a reduction in the number of household owners under the age of 35 , may lead to tens of thousands of pounds per year if today ’s <b>young consumers</b> report after a study of the HTA , as reported in the Financial Times , <b>the “ kernel age group for garden products ”</b> .

Table 8: In this sample translation, the proposed character-to-character model fails to adequately capture a long-term dependency.

mistakes (Table 7 (a)).

Given a long, rare word such as “Siebentausendzweihundertvierundfünfzig” (seven thousand two hundred fifty four) in Table 7 (b), the subword-level model segments “Siebentausend” as (Sieb, ent, aus, end), which results in an inaccurate translation. The character-level model performs better on these long, concatenative words with ambiguous segmentation.

We expect a character-level model to handle novel and unseen morphological inflections well. We observe that this is indeed the case, as our char2char model correctly understands “gesperrt”, a past participle form of “sperrern” (to block) (Table 7 (c)).

Nonce words are terms coined for a single use. They are not actual words but are constructed in a way that humans can intuitively guess what they mean, such as *workoliday* and *friyay*. We construct a few DE-EN sentence pairs that contain German nonce words (one example shown in Table 7 (d)), and observe that the character-level model can indeed detect salient character patterns and arrive at a correct translation.

Finally, we evaluate our multilingual models’ capacity to perform intra-sentence code-switching, by giving them as input mixed sentences from multiple languages. The newstest-2013 development datasets for DE-EN, CS-EN and FI-EN contain intersecting examples with the same English sentences. We compile a list of these sentences in DE/CS/FI and their translation in EN, and uniformly choose a few samples at random from the English side. Words or clauses from different languages are manually inter-

mixed to create multilingual sentences.

We discover that when given sentences with a high degree of language intermixing, as in Table 7 (e), the multilingual bpe2char model fails to seamlessly handle alternation of languages. Overall, however, both multilingual models generate reasonable translations. This is possible because we did not provide a language identifier when training our multilingual models. As a result, they learned to understand a multilingual sentence and translate it into a coherent English sentence.

There are indeed cases where the proposed character-level model fails, and we notice that those are often sentences with long-distance dependencies (see Table 8).

We show supplementary, sample translations in each scenario on a webpage.<sup>4</sup>

**Training and decoding speed** On a single Titan X GPU, we observe that our char2char models are approximately 35% slower to train than our bpe2char baselines when the same batch size was used. Our bilingual character-level models can be trained in roughly two weeks.

We further note that the bilingual bpe2char model can translate 3,000 sentences in 66.63 minutes while the bilingual char2char model requires 71.71 minutes (online, not in batch). See Table 9 for the exact details.

**Further observations** We also note that the mul-

<sup>4</sup><https://sites.google.com/site/dl4mtc2c>

	Model	Time to execute 1k updates (s)	Batch size	Time to decode 3k sentences (m)
FI-EN	bpe2char	2461.72	128	66.63
	char2char	2371.93	64	71.71
Multi	bpe2char	1646.37	64	68.99
	char2char	2514.23	64	72.33

Table 9: Speed comparison. The second column shows the time taken to execute 1,000 training updates. The model makes each update after having seen one mini-batch.

tiling models are less prone to overfitting than the bilingual models. This is particularly visible for low-resource language pairs such as FI-EN. Figure 2 shows the evolution of the FI-EN validation BLEU scores where the bilingual models overfit rapidly but the multilingual models seem to regularize learning by training simultaneously on other language pairs.

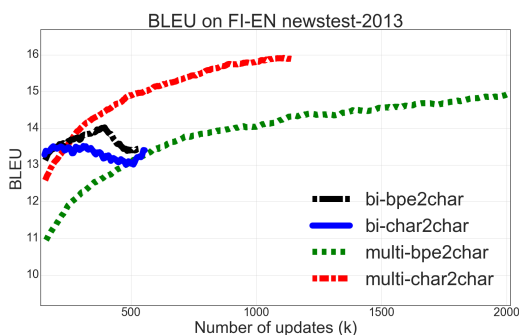


Figure 2: Multilingual models overfit less than bilingual models on low-resource language pairs.

## 8 Conclusion

We propose a fully character-level NMT model that accepts a sequence of characters in the source language and outputs a sequence of characters in the target language. What is remarkable about this model is the absence of explicitly hard-coded knowledge of words and their boundaries, and that the model learns these concepts from a translation task alone.

Our empirical results show that the fully character-level model performs as well as, or better than, subword-level translation models. The performance gain is distinctly pronounced in the multilingual many-to-one translation task, where results show that the character-level model can assign

model capacities to different languages more efficiently than the subword-level models. We observe a particularly large improvement in FI-EN translation when the model is trained to translate multiple languages, indicating a positive cross-lingual transfer to a low-resource language pair.

We discover two main benefits of the multilingual character-level model: (1) it is much more parameter-efficient than the bilingual models; and (2) it can naturally handle intra-sentence code-switching as a result of the many-to-one translation task. Ultimately, we present a case for fully character-level translation: that translation at the level of character is strongly beneficial and should be encouraged more.

The repository <https://github.com/nyu-dl/dl4mt-c2c> contains the source code and pre-trained models for reproducing the experimental results.

In the next stage of this research, we will investigate extending our multilingual many-to-one translation models to perform many-to-many translations, which will allow the decoder, similarly with the encoder, to learn from multiple target languages. Furthermore, a more thorough investigation into model architectures and hyperparameters is needed.

## Acknowledgements

Kyunghyun Cho thanks the support of eBay, Facebook, Google (Google Faculty Award, 2016) and NVidia (NVIDIA AI Lab, 2016-2019). This work was partly supported by the Samsung Advanced Institute of Technology (Deep Learning). Jason Lee was supported by the Qualcomm Innovation Fellowship, and thanks David Yenicelik and Kevin Wallmann for their contribution in designing the qualitative analysis. The authors would like to also thank Prof. Zheng Zhang (NYU, Shanghai) for fruitful discussions and comments, as well as Yvette Graham for her help with the human evaluation. Finally, the authors thank the Action Editor and anonymous reviewers for their constructive feedback.

## References

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *Proceedings of*

- the International Conference on Learning Representations.*
- Chris Callison-Burch. 2009. Fast, cheap, and creative: Evaluating translation quality using amazon’s mechanical turk. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 286–295.
- Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014a. On the properties of neural machine translation: Encoder-decoder approaches. In *Proceedings of the 8th Workshop on Syntax, Semantics, and Structure in Statistical Translation*, pages 103–111.
- Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014b. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of the Empirical Methods in Natural Language Processing*, pages 1724–1734.
- Junyoung Chung, Kyunghyun Cho, and Yoshua Bengio. 2016. A character-level decoder without explicit segmentation for neural machine translation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, pages 1693–1703.
- Marta R. Costa-Jussá and Josè A. R. Fonollosa. 2016. Character-based neural machine translation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, pages 357–361.
- Ferdinand de Saussure. 1916. *Course in General Linguistics*.
- Orhan Firat, Kyunghyun Cho, and Yoshua Bengio. 2016a. Multi-way, multilingual neural machine translation with a shared attention mechanism. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics*, pages 866–875.
- Orhan Firat, Baskaran Sankaran, Yaser Al-Onaizan, Fatos T. Yarman Vural, and Kyunghyun Cho. 2016b. Zero-resource translation with multi-lingual neural machine translation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 268–277.
- Dan Gillick, Cliff Brunk, Oriol Vinyals, and Amarnag Subramanya. 2015. Multilingual language processing from bytes. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics*, pages 1296–1306.
- Yvette Graham, Nitika Mathur, and Timothy Baldwin. 2015. Accurate evaluation of segment-level machine translation metrics. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1183–1191.
- Yvette Graham, Timothy Baldwin, Alistair Moffat, and Justin Zobel. 2017. Can machine translation systems be evaluated by the crowd alone. *Natural Language Engineering*, 23(1):3–30.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural Computation*, 9(8):1735–1780.
- Ray S. Jackendoff. 1992. *Semantic Structures*, Volume 18. MIT press.
- Sébastien Jean, Kyunghyun Cho, Roland Memisevic, and Yoshua Bengio. 2015. On using very large target vocabulary for neural machine translation. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics*, pages 1–5.
- Yoon Kim, Yacine Jernite, David Sontag, and Alexander M. Rush. 2015. Character-aware neural language models. In *Proceedings of the 30th AAAI Conference on Artificial Intelligence*, pages 2741–2749.
- Diederik Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference for Learning Representations*.
- Wang Ling, Isabel Trancoso, Chris Dyer, and Alan W. Black. 2015. Character-based neural machine translation. *arXiv preprint arXiv:1511.04586*.
- Minh-Thang Luong and Christopher D. Manning. 2016. Achieving open vocabulary neural machine translation with hybrid word-character models. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, pages 1054–1063.
- Minh-Thang Luong, Hieu Pham, and Christopher D. Manning. 2015. Effective approaches to attention-based neural machine translation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, pages 1412–1421.
- Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. 2013. On the difficulty of training recurrent neural networks. In *Proceedings of the 30th International Conference on Machine Learning*, pages 1310–1318.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, pages 1715–1725.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Edinburgh neural machine translation systems for WMT 16. In *Proceedings of the 1st Conference on Machine Translation*.
- Rupesh Kumar Srivastava, Klaus Greff, and Jürgen Schmidhuber. 2015. Training very deep networks. In *Advances in Neural Information Processing Systems*, Volume 28, pages 2377–2385.

- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, Volume 27, pages 3104–3112.
- Yulia Tsvetkov, Sunayana Sitaram, Manaal Faruqui, Guillaume Lample, Patrick Littell, David Mortensen, Alan W. Black, Lori Levin, and Chris Dyer. 2016. Polyglot neural language models: A case study in cross-lingual phonetic representation learning. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics*, pages 1357–1366.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.
- Yijun Xiao and Kyunghyun Cho. 2016. Efficient character-level document classification by combining convolution and recurrent layers. *arXiv preprint arXiv:1602.00367*.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In *Advances in Neural Information Processing Systems*, Volume 28, pages 649–657.