

Neural Machine Translation: A Review

Felix Stahlberg¹

University of Cambridge, Engineering Department, UK

Abstract

The field of machine translation (MT), the automatic translation of written text from one natural language into another, has experienced a major paradigm shift in recent years. Statistical MT, which mainly relies on various count-based models and which used to dominate MT research for decades, has largely been superseded by neural machine translation (NMT), which tackles translation with a single neural network. In this work we will trace back the origins of modern NMT architectures to word and sentence embeddings and earlier examples of the encoder-decoder network family. We will conclude with a survey of recent trends in the field.

Keywords: Neural machine translation, Neural sequence models

Various fields in the area of natural language processing (NLP) have been boosted by the rediscovery of neural networks [1]. However, for a long time, the integration of neural nets into machine translation (MT) systems was rather shallow. Early attempts used feedforward neural language models [2, 3] for the target language to rerank translation lattices [4]. The first neural models which also took the source language into account extended this idea by using the same model with bilingual tuples instead of target language words [5], scoring phrase pairs directly with a feedforward net [6], or adding a source context window to the neural language model [7, 8]. Kalchbrenner and Blunsom [9] and Cho et al. [10] introduced recurrent networks for translation modelling. All those approaches applied neural networks as component in a traditional statistical machine translation system. Therefore, they retained the log-linear model combination and only exchanged parts in the traditional architecture.

Neural machine translation (NMT) has overcome this separation by using a single large neural net that directly transforms the source sentence into the target sentence [11–13]. The advent of NMT certainly marks one of the major milestones in the history of MT, and has led to a radical and sudden departure of mainstream research from many previous research lines. This is perhaps best reflected by the explosion of scientific publications related to NMT in the past years² (Fig. 1), and the large number of publicly available NMT toolkits (Tab. 1). NMT has already been widely adopted in the industry [14–17] and is deployed in production systems by Google, Microsoft, Facebook, Amazon, SDL, Yandex, and many more. This article will introduce the basic concepts of NMT, and

¹Now at Google Research.

²Example Google Scholar search: https://scholar.google.com/scholar?q=%22neural+machine+translation%22&as_ylo=2017&as_yhi=2017

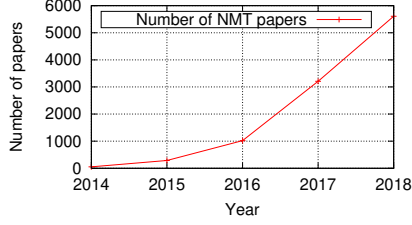


Figure 1: Number of papers mentioning “neural machine translation” per year according Google Scholar.

Name	Citation	Framework	GitHub Stars
Tensor2Tensor	Vaswani et al. [22]	TensorFlow	██████████
TensorFlow/NMT	-	TensorFlow	██████
Fairseq	Ott et al. [23]	PyTorch	██████
OpenNMT-py	Klein et al. [24]	Lua, (Py)Torch, TF	██████
Sockeye	Hieber et al. [25]	MXNet	██
OpenSeq2Seq	Kuchaiev et al. [26]	TensorFlow	██
Nematus	Sennrich et al. [27]	TensorFlow, Theano	██
PyTorch/Translate	-	PyTorch	██
Marian	Junczys-Dowmunt et al. [28]	C++	██
NMT-Keras	Álvaro Peris and Casacuberta [29]	TensorFlow, Theano	██
Neural Monkey	Helcl and Libovický [30]	TensorFlow	██
THUMT	Zhang et al. [31]	TensorFlow, Theano	██
Eske/Seq2Seq	-	TensorFlow	██
XNMT	Neubig et al. [32]	DyNet	██
NJUNMT	-	PyTorch, TensorFlow	██
Transformer-DyNet	-	DyNet	██
SGNMT	Stahlberg et al. [33, 34]	TensorFlow, Theano	██
CythonMT	Wang et al. [35]	C++	██
Neutron	Xu and Liu [36]	PyTorch	██

Table 1: NMT tools that have been updated in the past year (as of 2019). GitHub stars indicate the popularity of tools on GitHub.

will give a comprehensive overview of current research in the field. For even more insight into the field of neural machine translation, we refer the reader to other overview papers such as [18–21].

1. Nomenclature

We will denote the source sentence of length I as \mathbf{x} . We use the subscript i to index tokens in the source sentence. We refer to the source language vocabulary as Σ_{src} .

$$\mathbf{x} = x_1^I = (x_1, \dots, x_I) \in \Sigma_{src}^I \quad (1)$$

The translation of source sentence \mathbf{x} into the target language is denoted as \mathbf{y} . We use an analogous nomenclature on the target side.

$$\mathbf{y} = y_1^J = (y_1, \dots, y_J) \in \Sigma_{trg}^J \quad (2)$$

In case we deal with only one language we drop the subscript *src/trg*. For convenience we represent tokens as indices in a list of subwords or word surface forms. Therefore, Σ_{src} and Σ_{trg} are the first n natural numbers (i.e. $\Sigma = \{n' \in \mathbb{N} | n' \leq n\}$ where $n = |\Sigma|$ is the vocabulary size). Additionally, we use the projection function π_k which maps a tuple or vector to its k -th entry:

$$\pi_k(z_1, \dots, z_k, \dots, z_n) = z_k. \quad (3)$$

For a matrix $A \in \mathbb{R}^{m \times n}$ we denote the element in the p -th row and the q -th column as $A_{p,q}$, the p -th row vector as $A_{p,:} \in \mathbb{R}^n$ and the q -th column vector as $A_{:,q} \in \mathbb{R}^m$. For a series of m n -dimensional vectors $a_p \in \mathbb{R}^n$ ($p \in [1, m]$) we denote the $m \times n$ matrix which results from stacking the vectors horizontally as $(a_p)_{p=1:m}$ as illustrated with the following tautology:

$$A = (A_{p,:})_{p=1:m} = ((A_{:,q})_{q=1:n})^T. \quad (4)$$

2. Word Embeddings

Representing words or phrases as continuous vectors is arguably one of the keys in connectionist models for NLP. To the best of our knowledge, continuous space word representations were first successfully used for language modelling [2, 37]. The key idea is to represent a word $x \in \Sigma$ as a d -dimensional vector of real numbers. The size d of the embedding layer is normally chosen to be much smaller than the vocabulary size ($d \ll |\Sigma|$) in order to obtain interesting representations. The mapping from the word to its distributed representation can be represented by an embedding matrix $E \in \mathbb{R}^{d \times |\Sigma|}$ [38]. The x^{th} column of E (denoted as E_x) holds the d -dimensional representation for the word x .

Learned continuous word representations have the potential of capturing morphological, syntactic and semantic similarity across words [38]. In neural machine translation, embedding matrices are usually trained jointly with the rest of the network using back-propagation [39] and a gradient based optimizer such as stochastic gradient descent. In other areas of NLP, pre-trained word embeddings trained on unlabelled text have become ubiquitous [40]. Methods for training word embeddings on raw text often take the context into account in which the word occurs frequently [41, 42], or use cross-lingual information to improve embeddings [43, 44].

A newly emerging type of *contextualized* word embeddings [45, 46] is gaining popularity in various fields of NLP. Contextualized representations do not only depend on the word itself but on the entire input sentence. Thus, they cannot be described by a single embedding matrix but are usually generated by neural sequence models which have been trained under a language model objective. Most approaches either use LSTM [45, 47] or Transformer architectures [48, 49] but differ in the way these architectures are used to compute the word representations. Contextualized word embeddings have advanced the state-of-the-art in several NLP benchmarks [47, 49, 50]. Goldberg [51] showed that contextualized embeddings are remarkably sensitive to syntax. Choi et al. [52] reported gains from contextualizing word embeddings in NMT using a bag of words.

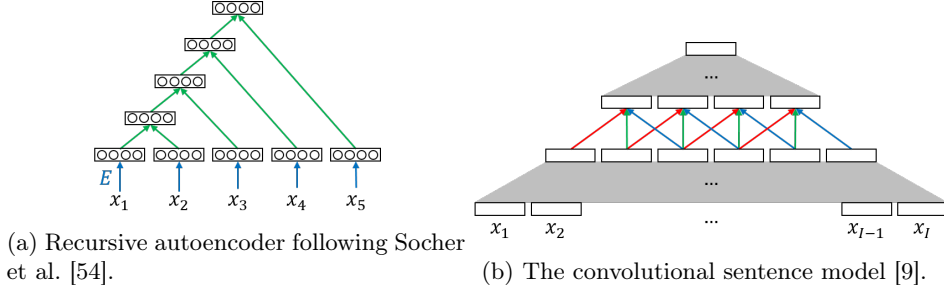


Figure 2: Phrase and sentence embedding architectures. The color coding indicates weight sharing.

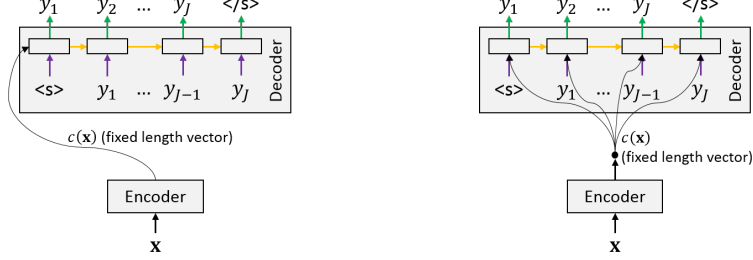
3. Phrase Embeddings

For various NLP tasks such as sentiment analysis or MT it is desirable to embed whole phrases or sentences instead of single words. For example, a distributed representation of the source sentence \mathbf{x} could be used as conditional for the distribution over the target sentences $P(\mathbf{y}|\mathbf{x})$. Early approaches to phrase embedding were based on recurrent autoencoders [53, 54]. To represent a phrase $\mathbf{x} \in \Sigma^I$ as d -dimensional vector, Socher et al. [54] first trained a word embedding matrix $E \in \mathbb{R}^{d \times |\Sigma|}$. Then, they recursively applied an autoencoder network which finds d -dimensional representations for $2d$ -dimensional inputs, where the input is the concatenation of two parent representations. The parent representations are either word embeddings or representations calculated by the same autoencoder from two different parents. The order in which representations are merged is determined by a binary tree over \mathbf{x} which can be constructed greedily [54] or derived from an Inversion Transduction Grammar [55, ITG] [56]. Fig. 2a shows an example of a recurrent autoencoder embedding a phrase with five words into a four dimensional space. One of the disadvantages of recurrent autoencoders is that the word and sentence embeddings need to have the same dimensionality. This restriction is not very critical in sentiment analysis because the sentence representation is only used to extract the sentiment of the writer [54]. In MT, however, the sentence representations need to convey enough information to condition the target sentence distribution on it, and thus should be higher dimensional than the word embeddings.

4. Sentence Embeddings

Kalchbrenner and Blunsom [9] used convolution to find vector representations of phrases or sentences and thus avoided the dimensionality issue of recurrent autoencoders. As shown in Fig. 2b, their model yields n -gram representations at each convolution level, with n increasing with depth. The top level can be used as representation for the whole sentence. Other notable examples of using convolution for sentence representations include [57–61]. However, the convolution operations in these models lose information about the exact word order, and are thus more suitable for sentiment analysis than for tasks like machine translation.³ A recent line of work uses self-attention rather than

³This is not to be confused with convolutional *translation* models which will be reviewed in Sec. 6.4



(a) Source sentence is used to initialize the decoder state. (b) Source sentence is fed to the decoder at each time step.

Figure 3: Encoder-decoder architectures with fixed-length sentence encodings. The color coding indicates weight sharing.

convolution to find sentence representations [62–64]. Another interesting idea explored by Yu et al. [65] is to resort to (recursive) relation networks [66, 67] which repeatedly aggregate pairwise relations between words in the sentence. Recurrent architectures are also commonly used for sentence representation. It has been noted that even random RNNs without any training can work surprisingly well for several NLP tasks [68–70].

5. Encoder-Decoder Networks with Fixed Length Sentence Encodings

Kalchbrenner and Blunsom [9] were the first who conditioned the target sentence distribution on a distributed fixed-length representation of the source sentence. Their recurrent continuous translation models (RCTM) I and II gave rise to a new family of so-called encoder-decoder networks which is the current prevailing architecture for NMT. Encoder-decoder networks are subdivided into an encoder network which computes a representation of the source sentence, and a decoder network which generates the target sentence from that representation. As introduced in Sec. 1 we denote the source sentence as $\mathbf{x} = x_1^J$ and the target sentence as $\mathbf{y} = y_1^J$. All existing NMT models define a probability distribution over the target sentences $P(\mathbf{y}|\mathbf{x})$ by factorizing it into conditionals:

$$P(\mathbf{y}|\mathbf{x}) \stackrel{\text{Chain rule}}{=} \prod_{j=1}^J P(y_j|y_1^{j-1}, \mathbf{x}). \quad (5)$$

Different encoder-decoder architectures differ vastly in how they model the distribution $P(y_j|y_1^{j-1}, \mathbf{x})$. We will first discuss encoder-decoder networks in which the encoder represents the source sentence as a fixed-length vector $c(\mathbf{x})$ like the methods in Sec. 4. The conditionals $P(y_j|y_1^{j-1}, \mathbf{x})$ are modelled as:

$$P(y_j|y_1^{j-1}, \mathbf{x}) = g(y_j|s_j, y_{j-1}, c(\mathbf{x})) \quad (6)$$

where s_j is the hidden state of a recurrent neural (decoder) network (RNN). We will formally introduce s_j in Sec. 6.3. Gated activation functions such as the long short-term memory [71, LSTM] or the gated recurrent unit [10, GRU] are commonly used to alleviate the vanishing gradient problem [72] which makes it difficult to train RNNs to

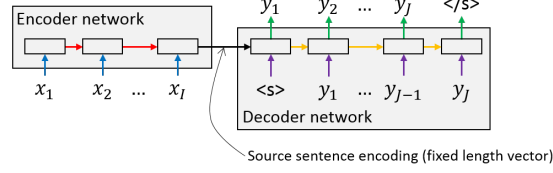
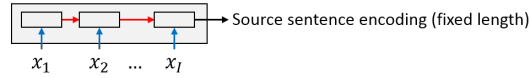
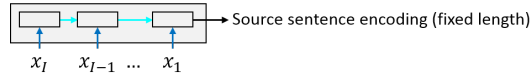


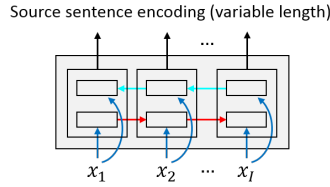
Figure 4: The encoder-decoder architecture of Sutskever et al. [12]. The color coding indicates weight sharing.



(a) Unidirectional encoder used by Cho et al. [10].



(b) Reversed unidirectional encoder from Sutskever et al. [12].



(c) Bidirectional encoder used by Bahdanau et al. [13].

Figure 5: Encoder architectures. The color coding indicates weight sharing.

capture long-range dependencies. Deep architectures with stacked LSTM cells were used by Sutskever et al. [12]. The encoder can be a convolutional network as in the RCTM I [9], an LSTM network [12], or a GRU network [10]. $g(\cdot)$ is a feedforward network with a softmax layer at the end which takes as input the decoder state s_j and an embedding of the previous target token y_{j-1} . In addition, $g(\cdot)$ may also take the source sentence encoding $c(\mathbf{x})$ as input to condition on the source sentence [9, 10]. Alternatively, $c(\mathbf{x})$ is just used to initialize the decoder state s_1 [12, 13]. Fig. 3 contrasts both methods. Intuitively, once the source sentence has been encoded, the decoder starts generating the first target sentence symbol y_1 which is then fed back to the decoder network for producing the second symbol y_2 . The algorithm terminates when the network produces the end-of-sentence symbol </s> . Sec. 7 explains more formally what we mean by the network “generating” a symbol y_j and sheds more light on the aspect of decoding in NMT. Fig. 4 shows the complete architecture of Sutskever et al. [12] who presented one of the first working standalone NMT systems that did not rely on any SMT baseline. One of the reasons why this paper was groundbreaking is the simplicity of the architecture, which stands in stark contrast to traditional SMT systems that used a very large number of highly engineered features.

Different ways of providing the source sentence to the encoder network have been explored in the past. Cho et al. [10] fed the tokens to the encoder in the natural order

they appear in the source sentence (cf. Fig. 5a). Sutskever et al. [12] reported gains from simply feeding the sequence in reversed order (cf. Fig. 5b). They argue that these improvements might be “caused by the introduction of many short term dependencies to the dataset” [12]. Bidirectional RNNs [73, BiRNN] are able to capture both directions (cf. Fig. 5c) and are often used in attentional NMT [13].

6. Attentional Encoder-Decoder Networks

6.1. Attention

One problem of early NMT models which is still not fully solved yet (see Sec. 10.1) is that they often produced poor translations for long sentences [74]. Cho et al. [11] suggested that this weakness is due to the fixed-length source sentence encoding. Sentences with varying length convey different amounts of information. Therefore, despite being appropriate for short sentences, a fixed-length vector “does not have enough capacity to encode a long sentence with complicated structure and meaning” [11]. Pouget-Abadie et al. [75] tried to mitigate this problem by chopping the source sentence into short clauses. They composed the target sentence by concatenating the separately translated clauses. However, this approach does not cope well with long-distance reorderings as word reorderings are only possible within a clause. Bahdanau et al. [13] introduced the concept of *attention* to avoid having a fixed-length source sentence representation. Their model does not use a constant context vector $c(\mathbf{x})$ any more which encodes the whole source sentence. By contrast, the attentional decoder can place its attention only on parts of the source sentence which are useful for producing the next token. The constant context vector $c(\mathbf{x})$ is thus replaced by a series of context vectors $c_j(\mathbf{x})$; one for each time step j .⁴

We will first introduce attention as a general concept before describing the architecture of Bahdanau et al. [13] in detail in Sec. 6.3. We follow the terminology of Vaswani et al. [76] and describe attention as mapping n query vectors to n output vectors via a mapping table (or a *memory*) of m key-value pairs. This view is related to memory-augmented neural networks which we will discuss in greater detail in Sec. 13.3. We make the simplifying assumption that all vectors have the same dimension d so that we can stack the vectors into matrices $Q \in \mathbb{R}^{n \times d}$, $K \in \mathbb{R}^{m \times d}$, and $V \in \mathbb{R}^{m \times d}$. Intuitively, for each query vector we compute an output vector as a weighted sum of the value vectors. The weights are determined by a similarity score between the query vector and the keys (cf. [76, Eq. 1]):

$$\underbrace{\text{Attention}(K, V, Q)}_{n \times d} = \text{Softmax}(\underbrace{\text{score}(Q, K)}_{n \times m}) \underbrace{V}_{m \times d}. \quad (7)$$

The output of $\text{score}(Q, K)$ is an $n \times m$ matrix of similarity scores. The softmax function normalizes over the columns of that matrix so that the weights for each query vector sum up to one. A straight-forward choice for $\text{score}(\cdot)$ proposed by Luong et al. [77] is the dot product (i.e. $\text{score}(Q, K) = QK^\top$). The most common scoring functions are summarized in Tab. 2.

⁴We refer to j as ‘time step’ due to the sequential structure of autoregressive models and the left-to-right order of NMT decoding. We note, however, that j does not specify a point in time in the usual sense but rather the position in the target sentence.

Name	Scoring function	Citation
Additive	$\text{score}(Q, K)_{p,q} = v^\top \tanh(WQ_{p,:} + UK_{q,:})$	Bahdanau et al. [13]
Dot-product	$\text{score}(Q, K) = QK^\top$	Luong et al. [77]
Scaled dot-product	$\text{score}(Q, K) = QK^\top d^{-0.5}$	Vaswani et al. [76]

Table 2: Common attention scoring functions. $v \in \mathbb{R}^{d_{\text{att}}}$, $W \in \mathbb{R}^{d_{\text{att}} \times d}$, and $U \in \mathbb{R}^{d_{\text{att}} \times d}$ in additive attention are trainable parameters with d_{att} being the dimensionality of the attention layer.

A common way to use attention in NMT is at the interface between encoder and decoder. Bahdanau et al. [13], Luong et al. [77] used the hidden decoder states s_j as query vectors. Both the key and value vectors are derived from the hidden states h_i of a recursive encoder.⁵ Formally, this means that $Q = s_j$ are the query vectors, $n = J$ is the target sentence length, $K = V = h_i$ are the key and value vectors, and $m = I$ is the source sentence length.⁶ The outputs of the attention layer are used as time-dependent context vectors $c_j(\mathbf{x})$. In other words, rather than using a fixed-length sentence encoding $c(\mathbf{x})$ as in Sec. 5, at each time step j we query a memory in which entries store (context-sensitive) representations of the source words. In this setup it is possible to derive an attention matrix $A \in \mathbb{R}^{J \times I}$ to visualize the learned relations between words in the source sentence and words in the target sentence:

$$A := \text{Softmax}(\text{score}((s_j)_{j=1:J}, (h_i)_{i=1:I})). \quad (8)$$

Fig. 6 shows an example of A from an English-German NMT system with additive attention. The attention matrix captures cross-lingual word relationships such as “is” \rightarrow “ist” or “great” \rightarrow “groß”. The system has learned that the English source word “is” is relevant for generating the German target word “ist” and thus emits a high attention weight for this pair. Consequently, the context vector $c_j(\mathbf{x})$ at time step $j = 3$ mainly represents the source word “is” ($c_3(\mathbf{x}) \approx h_2$). This is particularly significant as the system was not explicitly trained to align words but to optimize translation performance. However, as we will argue in Sec. 12.4, it would be wrong to think of A as a soft version of a traditional SMT word alignment.

An important generalization of attention is *multi-head* attention proposed by Vaswani et al. [76]. The idea is to perform H attention operations instead of a single one where H is the number of attention heads (usually $H = 8$). The query, key, and value vectors for the attention heads are linear transforms of Q , K , and V . The output of multi-head attention is the concatenation of the outputs of each attention head. The dimensionality of the attention heads is usually divided by H to avoid increasing the number of parameters. Formally, it can be described as follows [76]:

$$\text{MultiHeadAttention}(K, V, Q) = \text{Concat}(\text{head}_1, \dots, \text{head}_H)W^O \quad (9)$$

with weight matrix $W^O \in \mathbb{R}^{d \times d}$ where

$$\text{head}_h = \text{Attention}(KW_h^K, VW_h^V, QW_h^Q) \quad (10)$$

⁵ s_j and h_i are defined in Sec. 5 and Sec. 6.3.

⁶An exception is the model of Mino et al. [78] that splits h_i into two parts and uses the first part as key and the second as value.

	history	is	a	great	teacher	.	</s>
die							
Geschichte							
ist							
ein							
groß							
er							
Lehrer							
.							
</s>							

Figure 6: Attention weight matrix A for the translation from the English sentence “history is a great teacher .” to the German sentence “die Geschichte ist ein großer Lehrer .”. Dark shades of blue indicate high attention weights.

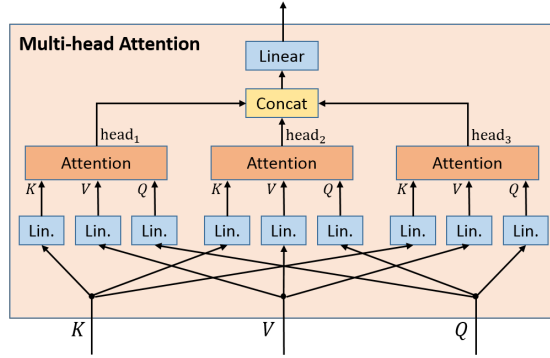


Figure 7: Multi-head attention with three attention heads.

with weight matrices $W_h^K, W_h^V, W_h^Q \in \mathbb{R}^{d \times \frac{d}{H}}$ for $h \in [1, H]$. Fig. 7 shows a multi-head attention module with three heads. Note that with multi-head attention it is not obvious anymore how to derive a single attention weight matrix A like shown in Fig. 6. Therefore, models using multi-head attention tend to be more difficult to interpret.

The concept of attention is no longer just a technique to improve sentence lengths in NMT. Since its introduction by Bahdanau et al. [13] it has become a vital part of various NMT architectures, culminating in the Transformer architecture (Sec. 6.5) which is entirely attention-based. Attention has also been proven effective for, inter alia, object recognition [79–81], image caption generation [82], video description [83], speech recognition [84, 85], cross-lingual word-to-phone alignment [86], bioinformatics [87], text summarization [88], text normalization [89], grammatical error correction [90], question answering [91–93], natural language understanding and inference [62, 94–96], uncertainty detection [97], photo optical character recognition [98], and natural language conversation [99].

6.2. Attention Masks and Padding

NMT usually groups sentences into batches to make more efficient use of the available hardware and to reduce noise in gradient estimation (cf. Sec. 11.1). However, the central data structure for many machine learning frameworks [101, 102] are *tensors* – multi-dimensional arrays with fixed dimensionality. Re-arranging source sentences as tensor

the	first	cold	shower	<pad>	<pad>
even	the	monkey	seems	to	want
a	little	coat	of	straw	<pad>

Figure 8: A tensor containing a batch of three source sentences of different lengths (“the first cold shower”, “even the monkey seems to want”, “a little coat of straw” – a haiku by Basho [100]). Short sentences are padded with <pad>. The training loss and attention masks are visualized with green (enabled) and red (disabled) background.

often results in some unused space as the sentences may vary in length. In practice, shorter sentences are filled up with a special padding symbol <pad> to match the length of the longest sentence in the batch (Fig. 8). Most implementations work with masks to avoid taking padded positions into account when computing the training loss. Attention layers also have to be restricted to non-padding symbols which is also usually realized by multiplying the attention weights by a mask that sets the attention weights for padding symbols to zero.

6.3. Recurrent Neural Machine Translation

This section contains a complete formal description of the RNNsearch architecture of Bahdanau et al. [13] which was the first NMT model using attention. Recall that NMT uses the chain rule to decompose the probability $P(\mathbf{y}|\mathbf{x})$ of a target sentence $\mathbf{y} = y_1^J$ given a source sentence $\mathbf{x} = x_1^I$ into left-to-right conditionals (Eq. 5). RNNsearch models the conditionals as follows [13, Eq. 2,4]:

$$P(\mathbf{y}|\mathbf{x}) \stackrel{\text{Eq. 5}}{=} \prod_{j=1}^J P(y_j|y_1^{j-1}, \mathbf{x}) = \prod_{j=1}^J g(y_j|y_{j-1}, s_j, c_j(\mathbf{x})). \quad (11)$$

Similarly to Eq. 6, the function $g(\cdot)$ encapsulates the decoder network which computes the distribution for the next target token y_j given the last produced token y_{j-1} , the RNN decoder state $s_j \in \mathbb{R}^n$, and the context vector $c_j(\mathbf{x}) \in \mathbb{R}^m$. The sizes of the encoder and decoder hidden layers are denoted with m and n . The context vector $c_j(\mathbf{x})$ is a distributed representation of the relevant parts of the source sentence. In NMT without attention [10, 12] (Sec. 5), the context vector is constant and thus needs to encode the whole source sentence. Adding an attention mechanism results in different context vectors for each target sentence position j . This effectively addresses issues in NMT due to the limited capacity of a fixed context vector as illustrated in Fig. 9.

As outlined in Sec. 6.1, the context vectors $c_j(\mathbf{x})$ are weighted sums of source sentence *annotations* $\mathbf{h} = (h_1, \dots, h_I)$. The annotations are produced by the encoder network. In other words, the encoder converts the input sequence \mathbf{x} to a sequence of annotations \mathbf{h} of the same length. Each annotation $h_i \in \mathbb{R}^m$ encodes information about the entire source sentence \mathbf{x} “with a strong focus on the parts surrounding the i -th word of the input sequence” [13, Sec. 3.1]. RNNsearch uses a bidirectional RNN [73, BiRNN] to generate the annotations. A BiRNN consists of two independent RNNs. The forward RNN \vec{f} reads \mathbf{x} in the original order (from x_1 to x_I). The backward RNN \overleftarrow{f} consumes \mathbf{x} in reversed order (from x_I to x_1):

$$\vec{h}_i = \vec{f}(x_i, \vec{h}_{i-1}) \quad (12)$$

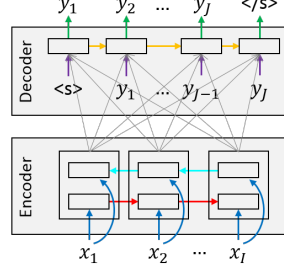


Figure 9: The RNNsearch model following Bahdanau et al. [13]. The color coding indicates weight sharing. Gray arrows represent attention.

$$\overleftarrow{h}_i = \overleftarrow{f}(x_i, \overleftarrow{h}_{i+1}). \quad (13)$$

The RNNs $\overrightarrow{f}(\cdot)$ and $\overleftarrow{f}(\cdot)$ are usually LSTM [71] or GRU [10] cells. The annotation h_i is the concatenation of the hidden states \overrightarrow{h}_i and \overleftarrow{h}_i [13, Sec. 3.2]:

$$h_i = [\overrightarrow{h}_i; \overleftarrow{h}_i]^\top. \quad (14)$$

The context vectors $c_j(\mathbf{x}) \in \mathbb{R}^m$ are computed from the annotations as weighted sum with weights $\alpha_j \in [0, 1]^I$ [13, Eq. 5]:

$$c_j(\mathbf{x}) = \sum_{i=1}^I \alpha_{j,i} h_i. \quad (15)$$

The weights are determined by the alignment model $a(\cdot)$:

$$\alpha_{j,i} = \frac{1}{Z} \exp(a(s_{j-1}, h_i)) \text{ with } Z = \sum_{k=1}^I \exp(a(s_{j-1}, h_k)) \quad (16)$$

where $a(s_{j-1}, h_i)$ is a feedforward neural network which estimates the importance of annotation h_i for producing the j -th target token given the current decoder state $s_{j-1} \in \mathbb{R}^n$. In the terminology of Sec. 6.1, h_i represent the keys and values, s_j are the queries, and $a(\cdot)$ is the attention scoring function.

The function $g(\cdot)$ in Eq. 11 does not only take the previous target token y_{j-1} and the context vector c_j but also the decoder hidden state s_j .

$$s_j = f(s_{j-1}, y_{j-1}, c_j) \quad (17)$$

where $f(\cdot)$ is modelled by a GRU or LSTM cell. The function $g(\cdot)$ is defined as follows.

$$g(y_j | y_{j-1}, s_j, c_j) \propto \exp(W_o \max(t_j, u_j)) \quad (18)$$

with

$$t_j = T_s s_j + T_y E y_{j-1} + T_c c_j \quad (19)$$

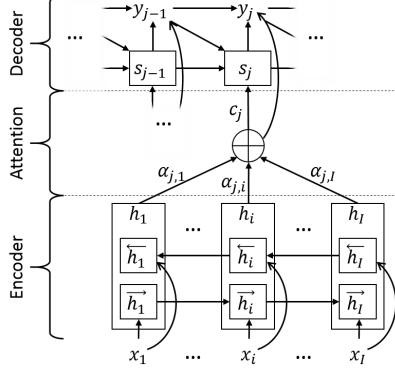


Figure 10: Illustration of the attention mechanism in RNNsearch [13].

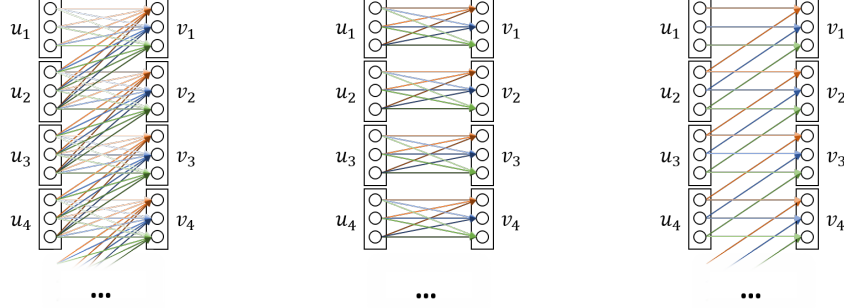
$$u_j = U_s s_j + U_y E y_{j-1} + U_c c_j \quad (20)$$

where $\max(\cdot)$ is the *element-wise* maximum, and $W_o \in \mathbb{R}^{|\Sigma_{trg}| \times l}$, $T_s, U_s \in \mathbb{R}^{l \times n}$, $T_y, U_y \in \mathbb{R}^{l \times k}$, $E \in \mathbb{R}^{k \times |\Sigma_{trg}|}$, $T_c, U_c \in \mathbb{R}^{l \times m}$ are weight matrices. The definition of $g(\cdot)$ can be seen as connecting the output of the recurrent layer, an k -dimensional embedding of the previous target token, and the context vector with a single maxout layer [103] of size l and using a softmax over the target language vocabulary [13]. Fig. 10 illustrates the complete RNNsearch model.

6.4. Convolutional Neural Machine Translation

Although convolutional neural networks (CNNs) have first been proposed by Waibel et al. [104] for phoneme recognition, their traditional use case is computer vision [105–107]. CNNs are especially useful for processing images because of two reasons. First, they use a high degree of weight tying and thus reduce the number of parameters dramatically compared to fully connected networks. This is crucial for high dimensional input like visual imagery. Second, they automatically learn space invariant features. Spatial invariance is desirable in vision since we often aim to recognize objects or features regardless of their exact position in the image. In NLP, convolutions are usually one dimensional since we are dealing with sequences rather than two dimensional images as in computer vision. We will therefore limit our discussions to the one dimensional case. We will also exclude concepts like pooling or strides as they are uncommon for sequence models in NLP.

The input to an 1D convolutional layer is a sequence of M -dimensional vectors u_1, \dots, u_I . The literature about CNNs usually refers to the M dimensions in each $u_i \in \mathbb{R}^M$ ($i \in [1, I]$) as *channels*, and to the i -axis as *spatial dimension*. The convolution transforms the input sequence u_1, \dots, u_I to an output sequence of N -dimensional v_1, \dots, v_I of the same length by moving a *kernel* of width K over the input sequence. The kernel is a linear transform which maps the K -gram u_i, \dots, u_{i+K-1} to the output v_i for $i \in [1, I]$ (we append $K - 1$ padding symbols to the input). Standard convolution



(a) Standard convolution. (b) Pointwise convolution. (c) Depthwise convolution.

Figure 11: Types of 1D-convolution used in NMT. The color coding indicates weight sharing.

Name	Number of parameters
Standard convolution	KMN
Pointwise convolution	MN
Depthwise convolution	KN
Depthwise separable convolution	N(M+K)

Table 3: Types of convolution and their number of parameters.

parameterizes this linear transform with a full weight matrix $W^{\text{std}} \in \mathbb{R}^{KM \times N}$:

$$\text{StdConv}:(v_i)_n = \sum_{m=1}^M \sum_{k=0}^{K-1} W_{kM+m,n}^{\text{std}} (u_{i+k})_m \quad (21)$$

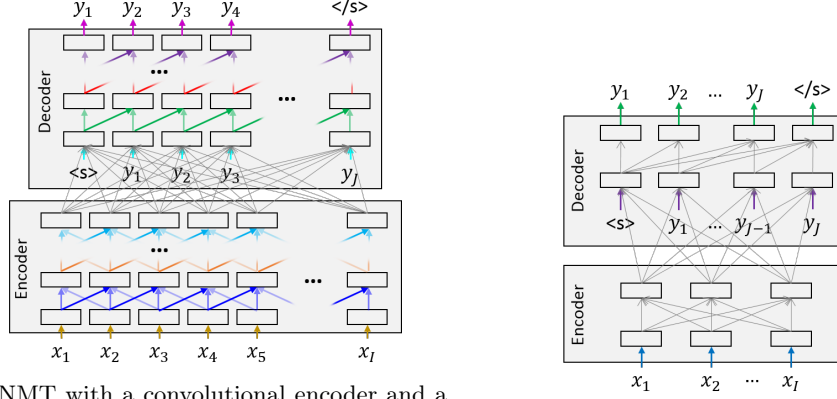
with $i \in [1, I]$ and $n \in [1, N]$. Standard convolution represents two kinds of dependencies: Spatial dependency (inner sum in Eq. 21) and cross-channel dependency (outer sum in Eq. 21). Pointwise and depthwise convolution factor out these dependencies into two separate operations:

$$\text{PointwiseConv}:(v_i)_n = \sum_{m=1}^M W_{m,n}^{\text{pw}} (u_i)_m = u_i W^{\text{pw}} \quad (22)$$

$$\text{DepthwiseConv}:(v_i)_n = \sum_{k=0}^{K-1} W_{k,n}^{\text{dw}} (u_{i+k})_n \quad (23)$$

where $W^{\text{pw}} \in \mathbb{R}^{M \times N}$ and $W^{\text{dw}} \in \mathbb{R}^{K \times N}$ are weight matrices. Fig. 11 illustrates the differences between these types of convolution. The idea behind *depthwise separable* convolution is to replace standard convolutional with depthwise convolution followed by pointwise convolution. As shown in Tab. 3, the decomposition into two simpler steps reduces the number of parameters and has been shown to make more efficient use of the parameters than regular convolution in vision [108, 109].

Using convolution rather than recurrence in NMT models has several potential advantages. First, they reduce sequential computation and are therefore easier parallelizable on



(a) NMT with a convolutional encoder and a convolutional decoder like in the ConvS2S architecture [110]. (b) Purely attention-based NMT as proposed by Vaswani et al. [76] with two layers.

Figure 12: Convolutional and purely attention-based architectures. The color coding indicates weight sharing. Gray arrows represent attention.

GPU hardware. Second, their hierarchical structure connects distant words via a shorter path than sequential topologies [110] which eases learning [72]. Both regular [110–112] and depthwise separable [113, 114] convolution have been used for NMT in the past. Fig. 12a shows the general architecture for a fully convolutional NMT model such as ConvS2S [110] or SliceNet [113] in which both encoder and decoder are convolutional. Stacking multiple convolutional layers increases the effective context size. In the decoder, we need to mask the receptive field of the convolution operations to make sure that the network has no access to future information [115]. Encoder and decoder are connected via attention. Gehring et al. [110] used attention into the encoder representations after each convolutional layer in the decoder.

6.5. Self-attention-based Neural Machine Translation

Recall that Eq. 5 states that NMT factorizes $P(\mathbf{y}|\mathbf{x})$ into conditionals $P(y_j|y_1^{j-1}, \mathbf{x})$. We have reviewed two ways to model the dependency on the source sentence \mathbf{x} in NMT: via a fixed-length sentence encoding $c(\mathbf{x})$ (Sec. 5) or via time-dependent context vectors $c_j(\mathbf{x})$ which are computed using attention (Sec. 6.1). We have also presented two ways to implement the dependency on the target sentence prefix y_1^{j-1} : via a recurrent connection which passes through the decoder state to the next time step (Sec. 6.3) or via convolution (Sec. 6.4). A third option to model target side dependency is using *self-attention*. Using the terminology introduced in Sec. 6.1, decoder self-attention derives all three components (queries, keys, and values) from the decoder state. **The decoder conditions on the translation prefix y_1^{j-1} by attending to its own states from previous time steps.** Besides machine translation, self-attention has been applied to various NLP tasks such as sentiment analysis [116], natural language inference [62, 96, 117, 118], text summarization [119], headline generation [120], sentence embedding [63, 64, 121], and reading comprehension [122]. **Similarly to convolution, self-attention introduces short paths between distant words and reduces the amount of sequential computation.** Studies

indicate that these short paths are especially useful for learning strong semantic feature extractors, but (perhaps somewhat counter-intuitively) less so for modelling long-range subject-verb agreement [123]. Like in convolutional models we also need to mask future decoder states to prevent conditioning on future tokens (cf. Sec. 6.2). The general layout for self-attention-based NMT models is shown in Fig. 12b. The first example of this new class of NMT models was the Transformer [76]. The Transformer uses attention for three purposes: 1) within the encoder to enable context-sensitive word representations which depend on the whole source sentence, 2) between the encoder and the decoder as in previous models, and 3) within the decoder to condition on the current translation history. The Transformer uses multi-head attention (Sec. 6.1) rather than regular attention. Using multi-head attention has been shown to be essential for the Transformer architecture [123, 124].

A challenge in self-attention-based models (and to some extent in convolutional models) is that vanilla attention as introduced in Sec. 6.1 by itself has no notion of order. The key-value pairs in the memory are accessed purely based on the correspondence between key and query (*content-based* addressing) and not based on a location of the key in the memory (*location-based*).⁷ This is less of a problem in recurrent NMT (Sec. 6.3) as queries, keys, and values are derived from RNN states and already carry a strong sequential signal due to the RNN topology. In the Transformer architecture, however, recurrent connections are removed in favor of attention. Vaswani et al. [76] tackled this problem using *positional encodings*. Positional encodings are (potentially partial) functions $PE : \mathbb{N} \rightarrow \mathbb{R}^D$ where D is the word embedding size, i.e. they are D -dimensional representations of natural numbers. They are added to the (input and output) word embeddings to make them (and consequently the queries, keys, and values) position-sensitive. Vaswani et al. [76] stacked sine and cosine functions of different frequencies to implement $PE(\cdot)$:

$$PE_{\sin}(n)_d = \begin{cases} \sin(10000^{-\frac{d}{D}} n) & : d \text{ is even} \\ \cos(10000^{-\frac{d}{D}} n) & : d \text{ is odd} \end{cases} \quad (24)$$

for $n \in \mathbb{N}$ and $d \in [1, D]$. Alternatively, positional encodings can be learned in an embedding matrix [110]:

$$PE_{\text{learned}}(n) = W_{:,n} \quad (25)$$

with weight matrix $W \in \mathbb{R}^{d \times N}$ for some sufficiently large N . The input to $PE(\cdot)$ is usually the absolute position of the word in the sentence [76, 110], but relative positioning is also possible [125]. We will give an overview of extensions to the Transformer architecture in Sec. 13.1.

6.6. Comparison of the Fundamental Architectures

As outlined in the previous sections, NMT can come in one of three flavors: recurrent, convolutional, or self-attention-based. In this section, we will discuss three concrete architectures in greater detail – one of each flavor. For an empirical comparison see [126]. Fig. 13 visualizes the data streams in Google’s Neural Machine Translation system [14,

⁷We will discuss cases in which both content and location are taken into account in Secs. 13.2 and 13.3

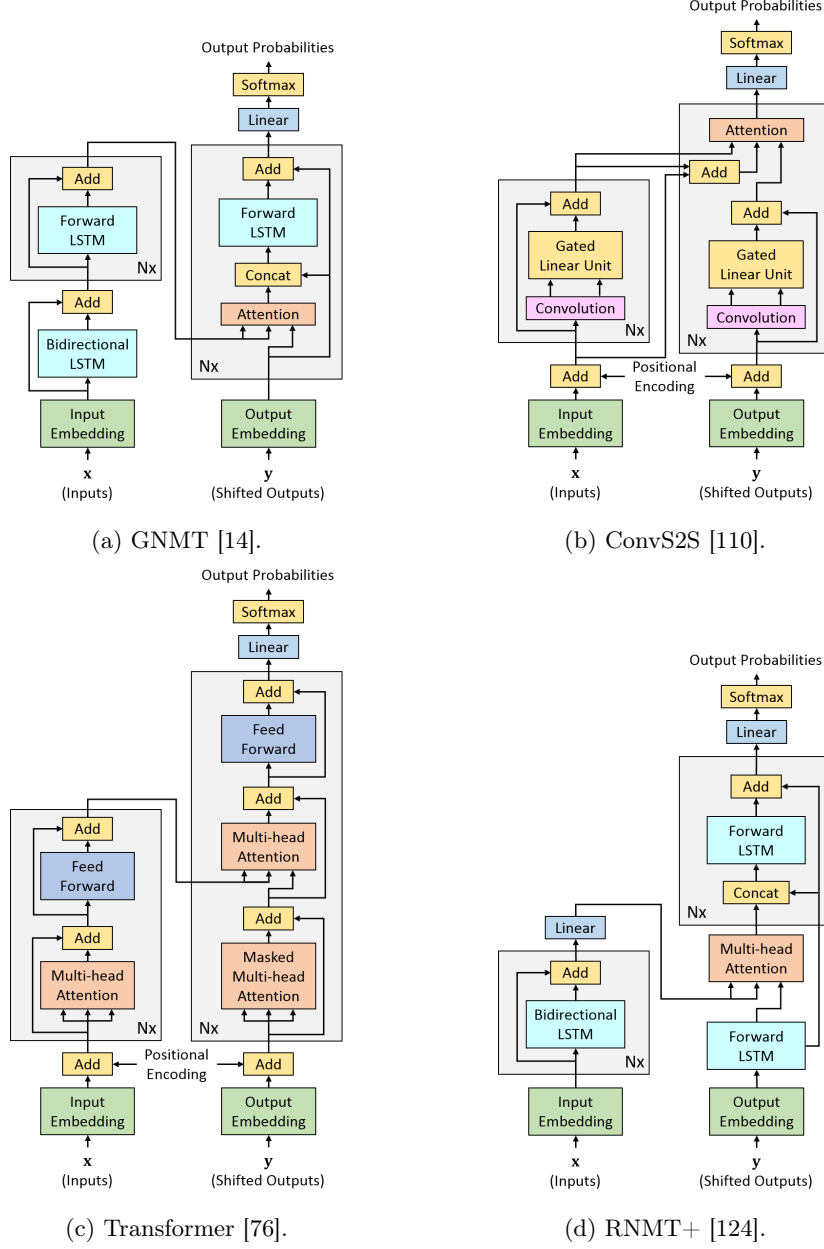


Figure 13: Comparison of NMT architectures. The three inputs to attention modules are (from left to right): keys (K), values (V), and queries (Q) as in Fig. 7.

GNMT] as example of a recurrent network, the convolutional ConvS2S model [110], and the self-attention-based Transformer model [76] in plate notation. We excluded components like dropout [127], batch normalization [128], and layer normalization [129] to

simplify the diagrams. All models fall in the general category of encoder-decoder networks, with the encoder in the left column and the decoder in the right column. Output probabilities are generated by a linear projection layer followed by a softmax activation at the end. They all use attention at each decoder layer to connect the encoder with the decoder, although the specifics differ. GNMT (Fig. 13a) uses regular attention, ConvS2S (Fig. 13b) adds the source word encodings to the values, and the Transformer (Fig. 13c) uses multi-head attention (Sec. 6.1). Residual connections [130] are used in all three architectures to encourage gradient flow in multi-layer networks. Positional encodings are used in ConvS2S and the Transformer, but not in GNMT. An interesting fusion is the RNMT+ model [124] shown in Fig. 13d which reintroduces ideas from the Transformer like multi-head attention into recurrent NMT. Other notable mixed architectures include Gehring et al. [112] who used a convolutional encoder with a recurrent decoder, Miculicich et al. [131], Wang et al. [132], Werlen et al. [133] who added self-attention connections to a recurrent decoder, Hao et al. [134] who used a Transformer encoder and a recurrent encoder in parallel, and Lin et al. [135] who equipped a recurrent decoder with a convolutional decoder to provide global target-side context.

7. Neural Machine Translation Decoding

7.1. The Search Problem in NMT

So far we have described how NMT defines the translation probability $P(\mathbf{y}|\mathbf{x})$. However, in order to apply these definitions directly, both the source sentence \mathbf{x} and the target sentence \mathbf{y} have to be given. They do not directly provide a method for generating a target sentence \mathbf{y} from a given source sentence \mathbf{x} which is the ultimate goal in machine translation. The task of finding the most likely translation $\hat{\mathbf{y}}$ for a given source sentence \mathbf{x} is known as the *decoding* or *inference* problem:

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y} \in \Sigma_{trg}^*} P(\mathbf{y}|\mathbf{x}). \quad (26)$$

NMT decoding is non-trivial for mainly two reasons. First, the search space is vast as it grows exponentially with the sequence length. For example, if we assume a common vocabulary size of $|\Sigma_{trg}| = 32,000$, there are already more possible translations with 20 words or less than atoms in the observable universe ($32,000^{20} \gg 10^{82}$). Thus, complete enumeration of the search space is impossible. Second, as we will see in Sec. 10, certain types of model errors are very common in NMT. The mismatch between the most likely and the “best” translation has deep implications on search as more exhaustive search often leads to worse translations [136]. We will discuss possible solutions to both problems in the remainder of Sec. 7.

7.2. Greedy and Beam Search

The most popular decoding algorithms for NMT are greedy search and beam search. Both search procedures are based on the left-to-right factorization of NMT in Eq. 5. Translations are built up from left to right while partial translation prefixes are scored using the conditionals $P(y_j|y_1^{j-1}, \mathbf{x})$. This means that both algorithms work in a time-synchronous manner: in each iteration j , partial hypotheses of (up to) length j are compared to each other, and a subset of them is selected for expansion in the next

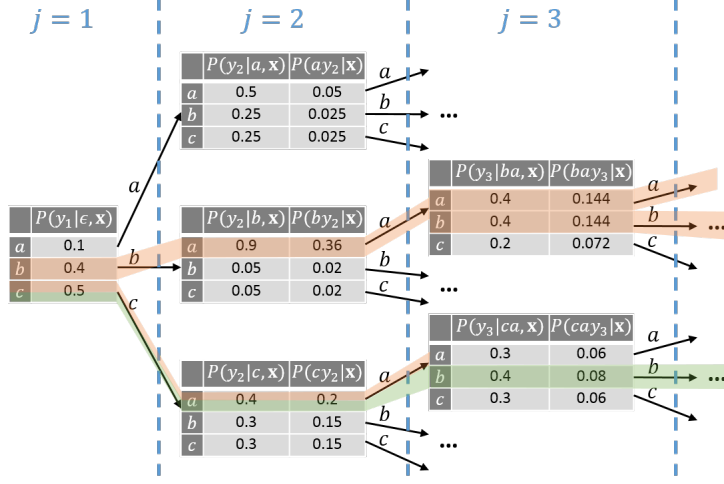


Figure 14: Comparison between greedy (highlighted in green) and beam search (highlighted in orange) with beam size 2.

time step. The algorithms terminate if either all or the best of the selected hypotheses end with the end-of-sentence symbol $\langle \text{</s>}$ or if some maximum number of iterations is reached. Fig. 14 illustrates the difference between greedy search and beam search. Greedy search (highlighted in green) selects the single best expansion at each time step: ‘c’ at $j = 1$, ‘a’ at $j = 2$, and ‘b’ at $j = 3$. However, greedy search is vulnerable to the so-called *garden-path problem* [20]. The algorithm selects ‘c’ in the first time step which turns out to be a mistake later on as subsequent distributions are very smooth and scores are comparably low. However, greedy decoding cannot correct this mistake later as it is already committed to this path. Beam search (highlighted in orange in Fig. 14) tries to mitigate the risk of the garden-path problem by passing not one but n possible translation prefixes to the next time step ($n = 2$ in Fig. 14). The n hypotheses which survive a time step are called *active hypotheses*. At each time step, the accumulated path scores for all possible continuations of active hypotheses are compared, and the n best ones are selected. Thus, beam search does not only expand ‘c’ but also ‘b’ in time step 1, and thereby finds the high scoring translation prefix ‘ba’. Note that although beam search seems to be the more accurate search procedure, it is not guaranteed to always find a translation with higher or equal score as greedy decoding.⁸ It is therefore still prone to the garden-path problem, although less so than greedy search. Stahlberg and Byrne [136] demonstrated that even beam search suffers from a high number of search errors.

7.3. Formal Description of Decoding for the RNNsearch Model

In this section, we will formally define decoding for the RNNsearch model [13]. We will resort to the mathematical symbols used in Sec. 6.3 to describe the algo-

⁸For example, imagine a series of high entropy conditionals after ‘baa’ and low entropy conditionals after ‘cab’ in Fig. 14

Algorithm 1 OneStepRNNsearch($s_{prev}, y_{prev}, \mathbf{h}$)

```
1:  $\alpha \xleftarrow{\text{Eq. 16}} \frac{1}{Z} [\exp(a(s_{prev}, h_i))]_{i \in [1, I]}$  {Attention weights ( $\alpha \in \mathbb{R}^I$ ,  $Z$  as in Eq. 16)}  
2:  $c \xleftarrow{\text{Eq. 15}} \sum_{i=1}^I \alpha_i \cdot h_i$  {Context vector update ( $c \in \mathbb{R}^m$ )}  
3:  $s \xleftarrow{\text{Eq. 17}} f(s_{prev}, y_{prev}, c)$  {RNN state update ( $s \in \mathbb{R}^n$ )}  
4:  $p \xleftarrow{\text{Eq. 5}} g(y_{prev}, s, c)$  { $p \in \mathbb{R}^{|\Sigma_{trg}|}$  is the distribution over the next target token  $P(y_j|\cdot)$ }  
5: return  $s, p$ 
```

Algorithm 2 GreedyRNNsearch(s_{init}, \mathbf{h})

```
1:  $\mathbf{y} \leftarrow \langle \rangle$   
2:  $s \leftarrow s_{init}$   
3:  $y \leftarrow \langle s \rangle$   
4: while  $y \neq \langle /s \rangle$  do  
5:    $s, p \leftarrow \text{OneStepRNNsearch}(s, y, \mathbf{h})$   
6:    $y \leftarrow \arg \max_{w \in \Sigma_{trg}} \pi_w(p)$   
7:    $\mathbf{y}.\text{append}(y)$   
8: end while  
9: return  $\mathbf{y}$ 
```

Algorithm 3 BeamRNNsearch($s_{init}, \mathbf{h}, n \in \mathbb{N}_+$)

```
1:  $\mathcal{H}_{cur} \leftarrow \{(\epsilon, 0.0, s_{init})\}$  {Initialize with empty translation prefix and zero score}  
2: repeat  
3:    $\mathcal{H}_{next} \leftarrow \emptyset$   
4:   for all  $(\mathbf{y}, p_{acc}, s) \in \mathcal{H}_{cur}$  do  
5:     if  $y_{|\mathbf{y}|} = \langle /s \rangle$  then  
6:        $\mathcal{H}_{next} \leftarrow \mathcal{H}_{next} \cup \{(\mathbf{y}, p_{acc}, s)\}$  {Hypotheses ending with  $\langle /s \rangle$  are not extended}  
7:     else  
8:        $s, p \leftarrow \text{OneStepRNNsearch}(s, y_{|\mathbf{y}|}, \mathbf{h})$   
9:        $\mathcal{H}_{next} \leftarrow \mathcal{H}_{next} \cup \bigcup_{w \in \Sigma_{trg}} (\mathbf{y} \cdot w, p_{acc} \pi_w(p), s)$  {Add all possible continuations}  
10:    end if  
11:  end for  
12:   $\mathcal{H}_{cur} \leftarrow \{(\mathbf{y}, p_{acc}, s) \in \mathcal{H}_{next} : |\{(\mathbf{y}', p'_{acc}, s') \in \mathcal{H}_{next} : p'_{acc} > p_{acc}\}| < n\}$  {Select  $n$ -best}  
13:   $(\hat{\mathbf{y}}, \hat{p}_{acc}, \hat{s}) \leftarrow \arg \max_{(\mathbf{y}, p_{acc}, s) \in \mathcal{H}_{cur}} p_{acc}$   
14: until  $\hat{y}_{|\hat{\mathbf{y}}|} = \langle /s \rangle$   
15: return  $\hat{\mathbf{y}}$ 
```

rithms. First, the source annotations \mathbf{h} are computed and stored as this does not require any search. Then, we compute the distribution for the first target token y_1 using OneStepRNNsearch($s_{init}, \langle s \rangle, \mathbf{h}$) (Alg. 1). The initial decoder state s_{init} is often a linear transform of the last encoder hidden state h_I : $s_{init} = Wh_I$ for some weight matrix $W \in \mathbb{R}^{n \times m}$.

Greedy decoding selects the most likely target token according the returned distribution and iteratively calls OneStepRNNsearch(\cdot) until the end-of-sentence symbol $\langle /s \rangle$ is emitted (Alg. 2). We use the projection function $\pi_w(p)$ (Eq. 3) which maps the posterior

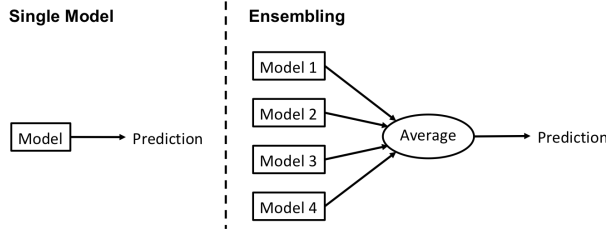


Figure 15: Ensembling four NMT models.

vector $p \in \mathbb{R}^{|\Sigma_{trg}|}$ to the w -th component.

The beam search strategy (Alg. 3) does not only keep the single best partial hypothesis but a set of n promising hypotheses where n is the size of the beam. A partial hypothesis is represented by a 3-tuple (\mathbf{y}, p_{acc}, s) with the translation prefix $\mathbf{y} \in \Sigma_{trg}^*$, the accumulated score $p_{acc} \in \mathbb{R}$, and the last decoder state $s \in \mathbb{R}^n$.

7.4. Ensembling

Ensembling [137, 138] is a simple yet very effective technique to improve the accuracy of NMT. The basic idea is illustrated in Fig. 15. The decoder makes use of K NMT networks rather than only one which are either trained independently [12, 14, 139] or share some amount of training iterations [140–142]. The ensemble decoder computes predictions for each of the individual models which are then combined using the arithmetic [12] or geometric [141] average:

$$S_{\text{arith}}(y_j | y_1^{j-1}, \mathbf{x}) = \frac{1}{K} \sum_{k=1}^K P_k(y_j | y_1^{j-1}, \mathbf{x}) \quad (27)$$

$$S_{\text{geo}}(y_j | y_1^{j-1}, \mathbf{x}) = \sum_{k=1}^K \log P_k(y_j | y_1^{j-1}, \mathbf{x}). \quad (28)$$

Both $S_{\text{arith}}(\cdot)$ and $S_{\text{geo}}(\cdot)$ can be used as drop-in replacement for the conditionals $P(y_j | y_1^{j-1}, \mathbf{x})$ in Eq. 5. The arithmetic average is more sound as $S_{\text{arith}}(\cdot)$ still forms a valid probability distribution which sums up to one. However, the geometric average $S_{\text{arith}}(\cdot)$ is numerically more stable as log-probabilities can be directly combined without converting them to probabilities. Note that the core idea of ensembling is similar to language model interpolation used in statistical machine translation or speech recognition.

Ensembling consistently outperforms single NMT by a large margin. All top systems in recent machine translation evaluation campaigns ensemble a number of NMT systems [126, 139–150], perhaps most famously taken to the extreme by the WMT18 submission of Tencent that ensembled up to 72 translation models [150]. However, the decoding speed is significantly worse since the decoder needs to apply K NMT models rather than only one. This means that the decoder has to perform K more forward passes through the networks, and has to apply the expensive softmax function K more times in each time step. Ensembling also often increases the number of CPU/GPU switches and the communication overhead between CPU and GPU when averaging is implemented on the CPU. Ensembling is also often more difficult to implement than single system NMT.

Knowledge distillation which we will discuss in Sec. 16 is one method to deal with the shortcomings of ensembling. Stahlberg and Byrne [151] proposed to unfold the ensemble into a single network and shrink the unfolded network afterwards for efficient ensembling.

In NMT, all models in an ensemble usually have the same size and topology and are trained on the same data. They differ only due to the random weight initialization and the randomized order of the training samples. Notable exceptions include Freitag and Al-Onaizan [152] who use ensembling to prevent overfitting in domain adaptation, He et al. [153] who combined models that selected their training data based on marginal likelihood, and the UCAM submission to WMT18 [126] that ensembled different NMT architectures with each other.⁹

When all models are equally powerful and are trained with the same data, it is surprising that ensembling is so effective. One common narrative is that different models make different mistakes, but the mistake of one model can be outvoted by the others in the ensemble [156]. This explanation is plausible for NMT since translation quality can vary widely between training runs [157]. The variance in translation performance may also indicate that the NMT error surface is highly non-convex such that the optimizer often ends up in local optima. Ensembling might mitigate this problem. Ensembling may also have a regularization effect on the final translation scores [158].

Checkpoint averaging [28, 159] is a technique which is often discussed in conjunction with ensembling [160]. Checkpoint averaging keeps track of the few most recent checkpoints during training, and averages their weight matrices to create the final model. This results in a single model and thus does not increase the decoding time. Therefore, it has become a very common technique in NMT [76, 126, 161]. Checkpoint averaging addresses a quite different problem than ensembling as it mainly smooths out minor fluctuations in the training curve which are due to the optimizer’s update rule or noise in the gradient estimation due to mini-batch training. In contrast, the weights of independently trained models are very different from each other, and there is no obvious direct correspondence between neuron activities across the models. Therefore, checkpoint averaging cannot be applied to independently trained models.

7.5. Decoding Direction

Standard NMT factorizes the probability $P(\mathbf{y}|\mathbf{x})$ from left to right (L2R) according Eq. 5. Mathematically, the left-to-right order is rather arbitrary, and other arrangements such as a right-to-left (R2L) factorization are equally correct:

$$\begin{aligned}
 P(\mathbf{y}|\mathbf{x}) &= \underbrace{\prod_{j=1}^J P(y_j|y_1^{j-1}, \mathbf{x})}_{=P(y_1|\mathbf{x}) \cdot P(y_2|y_1, \mathbf{x}) \cdot P(y_3|y_1, y_2, \mathbf{x}) \cdots} = \underbrace{\prod_{j=1}^J P(y_j|y_{j+1}^J, \mathbf{x})}_{=P(y_J|\mathbf{x}) \cdot P(y_{J-1}|y_J, \mathbf{x}) \cdot P(y_{J-2}|y_{J-1}, y_J, \mathbf{x}) \cdots} . \\
 &\hspace{15em} (29)
 \end{aligned}$$

NMT models which produce the target sentence in reverse order have led to some gains in evaluation systems when combined with left-to-right models [126, 140, 148, 150]. A common combination scheme is based on rescoring: A strong L2R ensemble first creates an n -best list which is then rescored with an R2L model [140, 162]. Stahlberg et al.

⁹Multi-source ensembling [154, 155] will be discussed in Sec. 15 in the context of multilingual NMT.

[126] used R2L models via a minimum Bayes risk framework. The L2R and R2L systems are normally trained independently, although some recent work proposes joint training schemes in which each direction is used as a regularizer for the other direction [163, 164]. Other orderings besides L2R and R2L have also been proposed such as middle-out [165], top-down in a binary tree [166], insertion-based [167–170], or in source sentence order [171].

Another way to give the decoder access to the full target-side context is the two-stage approach of Li et al. [172] who first drafted a translation, and then employed a multi-source NMT system to generate the final translation from both the source sentence and the draft. Zhang et al. [173] proposed a similar scheme but generated the draft translations in reverse order. A similar two-pass approach was used by ElMaghraby and Rafea [174] to make Arabic MT more robust against domain shifts. Geng et al. [175] used reinforcement learning to choose the best number of decoding passes.

Besides explicit combination with an R2L model and multi-pass strategies, we are aware of following efforts to make the decoder more sensitive to the right-side target context: He et al. [176] used reinforcement learning to estimate the long-term value of a candidate. Lin et al. [135] provided global target sentence information to a recurrent decoder via a convolutional model. Hoang et al. [177] proposed a very appealing theoretical framework to relax the discrete NMT optimization problem into a continuous optimization problem which allows to include both decoding directions.

7.6. Efficiency

NMT decoding is very fast on GPU hardware and can reach up to 5000 words per second.¹⁰ However, GPUs are very expensive, and speeding up CPU decoding to the level of SMT remains more challenging. Therefore, how to improve the efficiency of neural sequence decoding algorithms is still an active research question. One bottleneck is the sequential left-to-right order of beam search which makes parallelization difficult. Stern et al. [178] suggested to compute multiple time steps in parallel and validate translation prefixes afterwards. Kaiser et al. [179] reduced the amount of sequential computation by learning a sequence of latent discrete variables which is shorter than the actual target sentence, and generating the final sentence from this latent representation in parallel. Di Gangi and Federico [180] sped up recurrent NMT by using a simplified architecture for recurrent units. Another line of research tries to reintroduce the idea of *hypothesis recombination* to neural models. This technique is used extensively in traditional SMT [181]. The idea is to keep only the better of two partial hypotheses if it is guaranteed that both will be scored equally in the future. For example, this is the case for n -gram language models if both hypotheses end with the same n -gram. The problem in neural sequence models is that they condition on the full translation history. Therefore, hypothesis recombination for neural sequence models does not insist on exact equivalence but cluster hypotheses based on the similarity between RNN states or the n -gram history [182, 183]. A similar idea was used by Lecorvé and Motlicek [184] to approximate RNNs with WFSTs which also requires mapping histories into equivalence classes.

It is also possible to speed up beam search by reducing the beam size. Wu et al. [14], Freitag and Al-Onaizan [185] suggested to use a variable beam size, using various

¹⁰<https://marian-nmt.github.io/features/>

heuristics to decide the beam size at each time step. Alternatively, the NMT model training can be tailored towards the decoding algorithm [186–189]. Wiseman and Rush [187] proposed a loss function for NMT training which penalizes when the reference falls off the beam during training. Kim and Rush [190] reported that knowledge distillation (discussed in Sec. 16) reduces the gap between greedy decoding and beam decoding significantly. Greedy decoding can also be improved by using a small actor network which modifies the hidden states in an already trained model [189, 191].

7.7. *Generating Diverse Translations*

An issue with using beam search is that the hypotheses found by the decoder are very similar to each other and often differ only by one or two words [192–194]. The lack of diversity is problematic for several reasons. First, natural language in general and translation in particular often come with a high level of ambiguity that is not represented well by non-diverse n -best lists. Second, it impedes user interaction as NMT is not able to provide the user with alternative translations if needed. Third, collecting statistics about the search space such as estimating the probabilities of n -grams for minimum Bayes-risk decoding [126, 195–199] or risk-based training (Sec. 11.5) is much less effective.

Cho [200] added noise to the activations in the hidden layer of the decoder network to produce alternative high scoring hypotheses. This is justified by the observation that small variations of a hidden configuration encode semantically similar context [201]. Li and Jurafsky [192], Li et al. [193] proposed a diversity promoting modification of the beam search objective function. They added an explicit penalization term to the NMT score based on a maximum mutual information criterion which penalizes hypotheses from the same parent node. Note that both extensions can be used together [200]. Vijayakumar et al. [202] suggested to partition the active hypotheses in groups, and use a dissimilarity term to ensure diversity between groups. Park et al. [203] found alternative translations by k -nearest neighbor search from the greedy translation in a translation memory.

7.8. *Simultaneous Translation*

Most of the research in MT assumes an offline scenario: a complete source sentence is to be translated to a complete target sentence. However, this basic assumption does not hold up for many real-life applications. For example, useful machine translation for parliamentary speeches and lectures [204, 205] or voice call services such as Skype [206] does not only have to produce good translations but also have to do so with very low latency [207]. To reduce the latency in such real-time speech-to-speech translation scenarios it is desirable to start translating before the full source sentence has been vocalized by the speaker. Most approaches frame simultaneous machine translation as source sentence segmentation problem. The source sentence is revealed one word at a time. After a certain number of words, the segmentation policy decides to translate the current partial source sentence prefix and commit to a translation prefix which may not be a complete translation of the partial source. This process is repeated until the full source sentence is available. The segmentation policy can be heuristic [208] or learned with reinforcement learning [209, 210]. The translation itself is usually carried out by a standard MT system which was trained on full sentences. This is sub-optimal for two reasons. First, using a system which was trained on full sentences to translate partial sentences is brittle due to the significant mismatch between training and testing time. Ma et al. [211] tried to tackle

Vocabulary size	Number of parameters		
	Embeddings	Rest	Total
30K	55.8M	27.9M	83.7M
50K	93.1M	27.9M	121.0M
150K	279.2M	27.9M	307.1M

Table 4: Number of parameters in the original RNNsearch model [13] as presented in Sec. 6.3 (1000 hidden units, 620-dimensional embeddings). The model size highly depends on the vocabulary size.

this problem by training NMT to generate the target sentence with a fixed maximum latency to the source sentence. Second, human simultaneous interpreters use sophisticated strategies to reduce the latency by changing the grammatical structure [212–214]. These strategies are neglected by a vanilla translation system. Unfortunately, training data from human simultaneous translators is rare [213] which makes it difficult to adapt MT to it.

8. Open Vocabulary Neural Machine Translation

8.1. Using Large Output Vocabularies

As discussed in Sec. 2, NMT and other neural NLP models use embedding matrices to represent words as real-valued vectors. Embedding matrices need to have a fixed shape to make joint training with the translation model possible, and thus can only be used with a fixed and pre-defined vocabulary. This has several major implications for NMT.

First, the size of the embedding matrices grows with the vocabulary size. As shown in Tab. 4, the embedding matrices make up most of the model parameters of a standard RNNsearch model. Increasing the vocabulary size inflates the model drastically. Large models require a small batch size because they take more space in the (GPU) memory, but reducing the batch size often leads to noisier gradients, slower training, and eventually worse model performance [161]. Furthermore, a large softmax output layer is computationally very expensive. In contrast, traditional (symbolic) MT systems can easily use very large vocabularies [181, 215–217]. Besides these practical issues, training embedding matrices for large vocabularies is also complicated by the long-tail distribution of words in a language. Zipf’s law [218] states that the frequency of any word and its rank in the frequency table are inversely proportional to each other. Fig. 16 shows that 843K of the 875K distinct words (96.5%) occur less than 100 times in an English text with 140M running words – that is less than 0.00007% of the entire text. It is difficult to train robust word embeddings for such rare words. Word-based NMT models address this issue by restricting the vocabulary to the n most frequent words, and replacing all other words by a special token UNK. A problem with that approach is that the UNK token may appear in the generated translation. In fact, limiting the vocabulary to the 30K most frequent words results in an out-of-vocabulary rate (OOV) of 2.9% on the training set (Fig. 16). That means an UNK token can be expected to occur every 35 words. In practice, the number of UNKs is usually even higher. One simple reason is that the test set OOV rate is often higher than on the training set because the distribution of words and phrases naturally varies across genre, corpora, and time. Another observation is that word-based NMT often prefers emitting UNK even if a more appropriate word

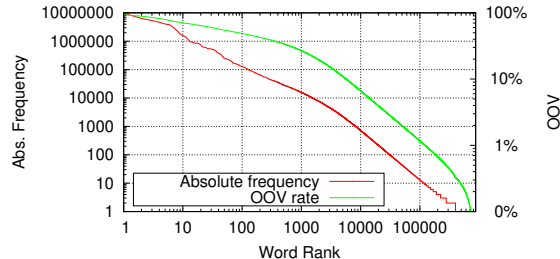


Figure 16: Distribution of words in the English portion of the English-German WMT18 training set (5.9M sentences, 140M words).

is in the NMT vocabulary. This is possibly due to the misbalance between the UNK token and other words: replacing all rare words with the same UNK token leads to an over-representation of UNK in the training set, and therefore a strong bias towards UNK during decoding.

8.1.1. Translation-specific Approaches

Jean et al. [219] distinguished between *translation-specific* and *model-specific* approaches. Translation-specific approaches keep the shortlist vocabulary in the original form, but correct UNK tokens afterwards. For example, the UNK replace technique [220, 221] keeps track of the positions of source sentence words which correspond to the UNK tokens. In a post-processing step, they replaced the UNK tokens with the most likely translation of the aligned source word according a bilingual word-level dictionary which was extracted from a word-aligned training corpus. Gulcehre et al. [222] followed a similar idea but used a special pointer network for referring to source sentence words. These approaches are rather ad-hoc because simple dictionary lookup without context is not a very strong model of translation. Li et al. [223] replaced each OOV word with a similar in-vocabulary word based on the cosine similarity between their distributed representations in a pre-processing step. However, this technique cannot tackle all OOVs as it is based on vector representations of words which are normally only available for a closed vocabulary. Moreover, the replacements might differ from the original meaning significantly. Further UNK replacement strategies were presented by Li et al. [224, 225], Miao et al. [226], but all share the inevitable limitation of all translation-specific approaches, namely that the translation model itself is indiscriminative between a large number of OOVs.

8.1.2. Model-specific Approaches

Model-specific approaches change the NMT model to make training with large vocabularies feasible. For example, Nguyen and Chiang [227] improved the translation of rare words in NMT by adding a lexical translation model which directly connects corresponding source and target words. Another very popular idea is to train networks to output probability distributions without using the full softmax [228]. Noise-contrastive estimation [229, 230, NCE] trains a logistic regression model which discriminates between real training examples and noise. For example, to train an embedding for a word w , Mnih and Kavukcuoglu [231] treat w as positive example, and sample from the global

unigram word distribution in the training data to generate negative examples. The logistic regression model is a binary classifier and thus does not need to sum over the full vocabulary. NCE has been used to train large vocabulary neural sequence models such as language models [232]. The technique falls into the category of self-normalizing training [228] because the model is trained to emit normalized distributions without explicitly summing over the output vocabulary. Self-normalization can also be achieved by adding the value of the partition function to the training loss [8], encouraging the network to learn parameters which generate normalized output.

Another approach (sometimes referred to as *vocabulary selection*) is to approximate the partition function of the full softmax by using only a subset of the vocabulary. This subset can be selected in different ways. For example, Jean et al. [219] applied importance sampling to select a small set of words for approximating the partition function. Both softmax sampling and UNK replace have been used in one of the winning systems at the WMT’15 evaluation on English-German [233]. Various methods have been proposed to select the vocabulary to normalize over during decoding, such as fetching all possible translations in a conventional phrase table [234], using the vocabulary of the translation lattices from a traditional MT system [235, local softmax], and attention-based [236] and embedding-based [237] methods.

8.2. Character-based NMT

Arguably, both translation-specific and model-specific approaches to word-based NMT are fundamentally flawed. Translation-specific techniques like UNK replace are indiscriminative between translations that differ only by OOV words. A translation model which assigns exactly the same score to a large number of hypotheses is of limited use by its own. Model-specific approaches suffer from the difficulty of training embeddings for rare words (Sec. 8.1). Compound or morpheme splitting [238, 239] can mitigate this issue only to a certain extent. More importantly, a fully-trained NMT system even with a very large vocabulary cannot be extended with new words. However, customizing systems to new domains (and thus new vocabularies) is a crucial requirement for commercial MT. Moreover, many OOV words are proper names which can be passed through untranslated. Hiero [217] and other symbolic systems can easily be extended with new words and phrases.

More recent attempts try to alleviate the vocabulary issue in NMT by departing from words as modelling units. These approaches decompose the word sequences into finer-grained units and model the translation between those instead of words. To the best of our knowledge, Ling et al. [240] were the first who proposed an NMT architecture which translates between sequences of characters. The core of their NMT network is still on the word-level, but the input and output embedding layers are replaced with subnetworks that compute word representations from the characters of the word. Such a subnetwork can be recurrent [240, 241] or convolutional [242, 243]. This idea was extended to a hybrid model by Luong and Manning [244] who used the standard lookup table embeddings for in-vocabulary words and the LSTM-based embeddings only for OOVs.

Having a word-level model at the core of a character-based system does circumvent the closed vocabulary restriction of purely word-based models, but it is still segmentation-dependent: The input text has to be preprocessed with a tokenizer that separates words by blank symbols in languages without word boundary markers, optionally applies compound or morpheme splitting in morphologically rich languages, and isolates punctuation

symbols. Since tokenization is by itself error-prone and can degrade the translation performance [245], it is desirable to design character-level systems that do not require any prior segmentation. Chung et al. [246] used a bi-scale recurrent neural network that is similar to dynamically segmenting the input using jointly learned gates between a slow and a fast recurrent layer. Lee et al. [247], Yang et al. [248] used convolution to achieve segmentation-free character-level NMT. Costa-jussà et al. [249] took character-level NMT one step further and used bytes rather than characters to help multilingual systems. Gulcehre et al. [250] added a planning mechanism to improve the attention weights between character-based encoders and decoders.

8.3. Subword-unit-based NMT

As compromise between characters and full words, compression methods like Huffman codes [251], word piece models [14, 252], or byte pair encoding [157, 253, BPE] can be used to transform the words to sequences of subword units. Subwords have been used rarely for traditional SMT [254–256], but are currently the most common translation units for NMT. **Byte pair encoding (BPE) initializes the set of available subword units with the character set of the language. This set is extended iteratively in subsequent merge operations. Each merge combines the two units with the highest number of co-occurrences in the text.**¹¹ **This process terminates when the desired vocabulary size is reached. This vocabulary size is often set empirically, but can also be tuned on data [258].**

Given a fixed BPE vocabulary, there are often multiple ways to segment an unseen text.¹² The ambiguity stems from the fact that symbols are still part of the vocabulary even after they are merged. Most BPE implementations select a segmentation greedily by preferring longer subword units. Interestingly, the ambiguity can also be used as source of noise for regularization. Kudo [259] reported surprisingly large gains by augmenting the training data with alternative subword segmentations and by decoding from multiple segmentations of the same source sentence.

Segmentation approaches differ in the level of constraints they impose on the subwords. A common constraint is that subwords cannot span over multiple words [157]. However, enforcing this constraint again requires a tokenizer which is a potential source of errors (see Sec. 8.2). **The SentencePiece model [260] is a tokenization-free subword model that is estimated on raw text. On the other side of the spectrum, it has been observed that automatically learned subwords generally do not correspond to linguistic entities such as morphemes, suffixes, affixes etc.** However, linguistically-motivated subword units [261–264] that also take morpheme boundaries into account do not always improve over completely data-driven ones.

8.4. Words, Subwords, or Characters?

There is no conclusive agreement in the literature whether characters or subwords are the better translation units for NMT. Tab. 5 summarizes some of the arguments. The tendency seems to be that character-based systems have the potential of outperforming

¹¹Wu and Zhao [257] proposed alternatives to the co-occurrence counts. The wordpiece model [14, 252] can also be seen as replacing the co-occurrence counts with a language model objective.

¹²This is not true for other subword compression algorithms. For example, Huffman codes [251] are prefix codes and thus unique.

Character-based NMT	Subword-based NMT
<ul style="list-style-type: none"> + Better at transliteration [265]. + Dynamic segmentation favors characters [266]. + More robust against noise [267, 268]. + Better modelling of morphology [267]. + Character-level decoders better than subword-based ones in some studies [246, 269]. – Character-based NMT computationally more expensive than subword-based NMT [269]. – More prone to vanishing gradients [246]. – Long-range dependencies have to be modelled over longer time-spans [247]. 	<ul style="list-style-type: none"> + More grammatical [265]. + Iterative BPE segmentation favors larger vocabulary sizes [258]. + Better at syntax [267]. + Tends to outperform character-based models in recent MT evaluations [143–145].

Table 5: Summary of studies comparing characters and subword-units for neural machine translation.

subword-based NMT, but they are technically difficult to deploy. Therefore, most systems in the WMT18 evaluation are based on subwords [145]. On a more profound level, we do see the shift towards small modelling units not without some concern. Chung et al. [246] noted that “we often have a priori belief that a word, or its segmented-out lexeme, is a basic unit of meaning, making it natural to approach translation as mapping from a sequence of source-language words to a sequence of target-language words.” Translation is the task of transferring *meaning* from one language to another, and it makes intuitive sense to model this process with meaningful units. The decades of research in traditional SMT were characterized by a constant movement towards larger translation units – starting from the word-based IBM models [270] to phrase-based MT [181] and hierarchical SMT [217] that models syntactic structures. Expressions consisting of multiple words are even more appropriate units than words for translation since there is rarely a 1:1 correspondence between source and target words. In contrast, the starting point for character- and subword-based models is the language’s writing system. Most writing systems are not logographic but alphabetic or syllabic and thus use symbols without any relation to meaning. The introduction of symbolic word-level and phrase-level information to NMT is one of the main motivations for NMT-SMT hybrid systems (Sec. 18).

9. Using Monolingual Training Data

In practice, parallel training data for MT is hard to acquire and expensive, whereas untranslated monolingual data is usually abundant. This is one of the reasons why language models (LMs) are central to traditional SMT. For example, in Hiero [217], the translation grammar spans a vast space of possible translations but is weak in assigning scores to them. The LM is mainly responsible for selecting a coherent and fluent

translation from that space. However, the vanilla NMT formalism does not allow the integration of an LM or monolingual data in general.

There are several lines of research which investigate the use of monolingual training data in NMT. Gulcehre et al. [271, 272] suggested to integrate a separately trained RNN-LM into the NMT decoder. Similarly to traditional SMT [181] they started out with combining RNN-LM and NMT scores via a log-linear model (‘shallow fusion’). They reported even better performance with ‘deep fusion’ which uses a controller network that dynamically adjusts the weights between RNN-LM and NMT. Both deep fusion and n -best reranking with count-based language models have led to some gains in WMT evaluation systems [148, 233]. The ‘simple fusion’ technique [273] trains the translation model to predict the residual probability of the training data added to the prediction of a pre-trained and fixed LM.

The second line of research makes use of monolingual text via data augmentation. The idea is to add monolingual data in the target language to the natural parallel training corpus. Different strategies for filling in the source side for these sentences have been proposed such as using a single dummy token [274] or copying the target sentence over to the source side [275]. The most successful strategy is called back-translation [274, 276] which employs a separate translation system in the reverse direction to generate the source sentences for the monolingual target language sentences. The back-translating system is usually smaller and computationally cheaper than the final system for practical reasons, although with enough computational resources improving the quality of the reverse system can affect the final translation performance significantly [277]. Iterative approaches that back-translate with systems that were by themselves trained with back-translation can yield improvements [278–280] although they are not widely used due to their computational costs. Back-translation has become a very common technique and has been used in nearly all neural submissions to recent evaluation campaigns [140, 144, 145].

A major limitation of back-translation is that the amount of synthetic data has to be balanced with the amount of real parallel data [140, 274, 281]. Therefore, the back-translation technique can only make use of a small fraction of the available monolingual data. A misbalance between synthetic and real data can be partially corrected by over-sampling – duplicating real training samples a number of times to match the synthetic data size. However, very high over-sampling rates often do not work well in practice. Recently, Edunov et al. [282] proposed to add noise to the back-translated sentences to provide a stronger training signal from the synthetic sentence pairs. They showed that adding noise does not only improve the translation quality but also makes the training more robust against a high ratio of synthetic against real sentences. The effectiveness of using noise for data augmentation in NMT has also been confirmed by Wang et al. [283]. These methods increase the variety of the training data and thus make it harder for the model to fit which ultimately leads to stronger training signals. The variety of synthetic sentences in back-translation can also be increased by sampling multiple sentences from the reverse translation model [284].

A third class of approaches changes the NMT training loss function to incorporate monolingual data. For example, Cheng et al. [285], Tu et al. [286], Escolano et al. [287] proposed to add autoencoder terms to the training objective which capture how well a sentence can be reconstructed from its translated representation. Using the reconstruction error is also central to (unsupervised) dual learning approaches [288–290].

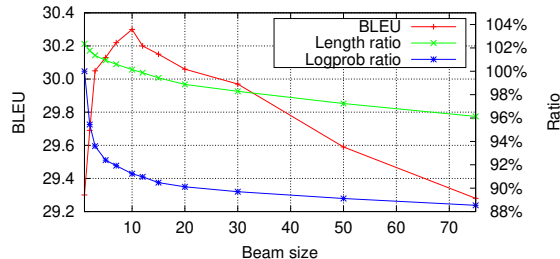


Figure 17: Performance of a Transformer model on English-German (WMT15) under varying beam sizes. The BLEU score peaks at beam size 10, but then suffers from a length ratio (hypothesis length / reference length) below 1. The log-probabilities are shown as a ratio with respect to greedy decoding.

However, training with respect to the new loss is often computationally intensive and requires approximations. Alternatively, multi-task learning has been used to incorporate source-side [291] and target-side [292] monolingual data. Another way of utilizing monolingual data in both source and target language is to warm start Seq2Seq training from pre-trained encoder and decoder networks [293, 294]. An extreme form of leveraging monolingual training data is unsupervised NMT which removes the need for parallel training data entirely. We will discuss unsupervised NMT in Sec. 14.4.

10. NMT Model Errors

NMT is highly effective in assigning scores (or probabilities) to translations because, in stark contrast to SMT, it does not make any conditional independence assumptions in Eq. 5 to model sentence-level translation.¹³ A potential drawback of such a powerful model is that it prohibits the use of sophisticated search procedures. Compared to hierarchical SMT systems like Hiero [217] that explore very large search spaces, NMT beam search appears to be overly simplistic. This observation suggests that translation errors in NMT are more likely due to *search errors* (the decoder does not find the highest scoring translation) than *model errors* (the model assigns a higher probability to a worse translation). Interestingly, this is not necessarily the case. Search errors in NMT have been studied by Stahlberg et al. [34], Stahlberg and Byrne [136], Niehues et al. [295]. In particular, Stahlberg and Byrne [136] demonstrated the high number of search errors in NMT decoding. However, as we will show in this section, NMT also suffers from various kinds of model errors in practice despite its theoretical advantage.

10.1. Sentence Length

Increasing the beam size exposes one of the most noticeable model errors in NMT. The red curve in Fig. 17 plots the BLEU score [296] of a recent Transformer NMT model against the beam size. A beam size of 10 is optimal on this test set. Wider beams lead to a steady drop in translation performance because the generated translations are becoming too short (green curve). However, as expected, the log-probabilities of the found

¹³It does, however, assume that each sentence can be translated in isolation. We will take a closer look at this assumption in Sec. 17.4.

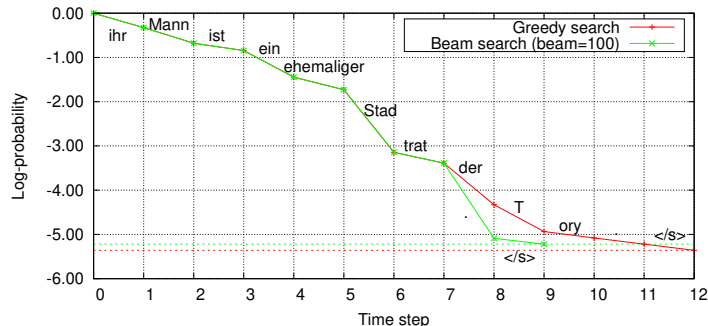


Figure 18: The length deficiency in NMT translating the English source sentence “Her husband is a former Tory councillor.” into German following Murray and Chiang [297]. The NMT model assigns a better score to the short translation “Ihr Mann ist ein ehemaliger Stadtrat.” than to the greedy translation “Ihr Mann ist ein ehemaliger Stadtrat der Tory.” even though it misses the former affiliation of the husband with the Tory Party.

translations (blue curve) are decreasing as we increase the beam size. NMT seems to assign too much probability mass to short hypotheses which are only found with more exhaustive search. Sountsov and Sarawagi [74] argue that this model error is due to the locally normalized maximum likelihood training objective in NMT that underestimates the margin between the correct translation and shorter ones if trained with regularization and finite data. A similar argument was made by Murray and Chiang [297] who pointed out the difficulty for a locally normalized model to estimate the “budget” for all remaining (longer) translations in each time step. Kumar and Sarawagi [298] demonstrated that NMT models are often poorly calibrated, and that calibration issues can cause the length deficiency in NMT. A similar case is illustrated in Fig. 18. The NMT model underestimates the combined probability mass of translations continuing after “Stadtrat” in time step 7 and overestimates the probability of the period symbol. Greedy decoding does not follow the green translation since “der” is more likely in time step 7. However, beam search with a large beam keeps the green path and thus finds the shorter (incomplete) translation with better score. In fact, Stahlberg and Byrne [136] linked the bias of large beam sizes towards short translations with the reduction of search errors.

At first glance this seems to be good news: fast beam search with a small beam size is already able to find good translations. However, fixing the model error of short translations by introducing search errors with a narrow beam seems like fighting fire with fire. In practice, this means that the beam size is yet another hyper-parameter which needs to be tuned for each new NMT training technique (eg. label smoothing [299] usually requires a larger beam), NMT architecture (the Transformer model is usually decoded with a smaller beam than typical recurrent models), and language pair [300]. More importantly, it is not clear whether there are gains to be had from reducing the number of search errors with wider beams which are simply obliterated by the NMT length deficiency.

10.1.1. Model-agnostic Length Models

The first class of approaches to alleviate the length problem is model-agnostic. Methods in this class treat the NMT model as black box but add a correction term to the NMT

score to bias beam search towards longer translations. A simple method is called *length normalization* which divides the NMT probability by the sentence length [233, 301]:

$$S_{\text{LN}}(\mathbf{y}|\mathbf{x}) = \frac{\log P(\mathbf{y}|\mathbf{x})}{|\mathbf{y}|} \quad (30)$$

Wu et al. [14] proposed an extension of this idea by introducing a tunable parameter α :

$$S_{\text{LN-GNMT}}(\mathbf{y}|\mathbf{x}) = \log P(\mathbf{y}|\mathbf{x}) \frac{(1 + 5)^\alpha}{(1 + |\mathbf{y}|)^\alpha} \quad (31)$$

Alternatively, like in SMT we can use a word penalty $\gamma(j, \mathbf{x})$ which rewards each word in the sentence:

$$S_{\text{WP}}(\mathbf{y}|\mathbf{x}) = \sum_{j=1}^J \gamma(j, \mathbf{x}) + \log P(y_j | y_1^{j-1}, \mathbf{x}) \quad (32)$$

A constant reward which is independent of \mathbf{x} and j can be found with the standard minimum-error-rate-training [302, MERT] algorithm [303] or with a gradient-based learning scheme [297]. Alternative policies which reward words with respect to some estimated sentence length were suggested by Huang et al. [304], Yang et al. [305].

10.1.2. Source-side Coverage Models

Tu et al. [306] connected the sentence length issue in NMT with the lack of an explicit mechanism to check the source-side coverage of a translation. Traditional SMT keeps track of a coverage vector $\mathcal{C}_{\text{SMT}} \in \{0, 1\}^I$ which contains 1 for source words which are already translated and 0 otherwise. \mathcal{C}_{SMT} is used to guard against *under-translation* (missing translations of some words) and *over-translation* (some words are unnecessarily translated multiple times). Since vanilla NMT does not use an explicit coverage vector it can be prone to both under- and over-translation [306, 307] and tends to prefer fluency over adequacy [308]. There are two popular ways to model coverage in NMT, both make use of the encoder-decoder attention weight matrix A introduced in Sec. 6.1. The simpler methods combine the scores of an already trained NMT system with a coverage penalty $cp(\mathbf{x}, \mathbf{y})$ without retraining. This penalty represents how much of the source sentence is already translated. Wu et al. [14] proposed the following term:

$$cp(\mathbf{x}, \mathbf{y}) = \beta \sum_{i=1}^I \log \left(\min \left(\sum_{j=1}^J A_{i,j}, 1.0 \right) \right). \quad (33)$$

A very similar penalty was suggested by Li et al. [309]:

$$cp(\mathbf{x}, \mathbf{y}) = \alpha \sum_{i=1}^I \log \left(\max \left(\sum_{j=1}^J A_{i,j}, \beta \right) \right) \quad (34)$$

where α and β are hyper-parameters that are tuned on the development set.

An even tighter integration can be achieved by changing the NMT architecture itself and jointly training it with a coverage model [306, 310]. Tu et al. [306] reintroduced an explicit coverage matrix $\mathcal{C} \in [0, 1]^{I \times J}$ to NMT. Intuitively, the j -th column $\mathcal{C}_{:,j}$ stores to

what extend each source word has been translated in time step j . \mathcal{C} can be filled with an RNN-based controller network (the “neural network based” coverage model of Tu et al. [306]). Alternatively, we can directly use A to compute the coverage (the “linguistic” coverage model of Tu et al. [306]):

$$c_{i,j} = \frac{1}{\Phi_i} \sum_{k=1}^j A_{i,k} \quad (35)$$

where Φ_i is the estimated number of target words the i -th source word generates which is similar to fertility in SMT. Φ_i is predicted by a feedforward network that conditions on the i -th encoder state. In both the neural network based and the linguistic coverage model, the decoder is modified to additionally condition on \mathcal{C} . The idea of using fertilities to prevent over- and under-translation has also been explored by Malaviya et al. [311]. A coverage model for character-based NMT was suggested by Kazimi and Costa-Jussá [312].

All approaches discussed in this section operate on the attention weight matrix A and are thus only readily applicable to models with single encoder-decoder attention like GNMT, but not to models with multiple encoder-decoder attention modules such as ConvS2S or the Transformer (see Sec. 6.6 for detailed descriptions of GNMT, ConvS2S, and the Transformer).

10.1.3. Controlling Mechanisms for Output Length

In some sequence prediction tasks such as headline generation or text summarization, the approximate desired output length is known in advance. In such cases, it is possible to control the length of the output sequence by explicitly feeding in the desired length to the neural model. The length information can be provided as additional input to the decoder network [313, 314], at each time step as the number of remaining tokens [315], or by modifying Transformer positional embeddings [316]. However, these approaches are not directly applicable to machine translation as the translation length is difficult to predict with sufficient accuracy.

11. NMT Training

NMT models are normally trained using backpropagation [39] and a gradient-based optimizer like Adadelta [317] with cross-entropy loss (Sec. 11.1). Modern NMT architectures like the Transformer, ConvS2S, or recurrent networks with LSTM [71] or GRU [10] cells help to address known training problems like vanishing gradients [72]. However, there is evidence that the optimizer still fails to exploit the full potential of NMT models and often gets stuck in suboptima:

1. NMT models vary greatly in performance, even if they use exactly the same architecture, training data, and are trained for the same number of iterations. Sennrich et al. [157] observed up to 1 BLEU difference between different models.
2. NMT ensembling (Sec. 15) combines the scores of multiple separately trained NMT models of the same kind. NMT ensembles consistently outperform single NMT by

a large margin. The achieved gains through ensembling might indicate difficulties in training of the single models.¹⁴

Training is therefore still a very active and diverse research topic. We will outline the different efforts in the literature on NMT training in this section.

11.1. Cross-entropy Training

The most common objective function for NMT training is cross-entropy loss. The optimization problem over model parameters Θ for a single sentence pair (\mathbf{x}, \mathbf{y}) under this loss is defined as follows:

$$\arg \min_{\Theta} \mathcal{L}_{\text{CE}}(\mathbf{x}, \mathbf{y}, \Theta) = \arg \min_{\Theta} - \sum_{j=1}^{|\mathbf{y}|} \log P_{\Theta}(y_j | y_1^{j-1}, \mathbf{x}). \quad (36)$$

In practice, NMT training groups several instances from the training corpus into batches, and optimizes Θ by following the gradient of the average $\mathcal{L}_{\text{CE}}(\mathbf{x}, \mathbf{y}, \Theta)$ in the batch. There are various ways to interpret this loss function.

Cross-entropy loss maximizes the log-likelihood of the training data. A direct interpretation of Eq. 36 is that it yields a maximum likelihood estimate of Θ as it directly maximizes the probability $P_{\Theta}(\mathbf{y}|\mathbf{x})$:

$$-\log P_{\Theta}(\mathbf{y}|\mathbf{x}) \stackrel{\text{Eq. 5}}{=} - \sum_{j=1}^{|\mathbf{y}|} \log P_{\Theta}(y_j | y_1^{j-1}, \mathbf{x}) = \mathcal{L}_{\text{CE}}(\mathbf{x}, \mathbf{y}, \Theta). \quad (37)$$

Cross-entropy loss optimizes a Monte Carlo approximation of the cross-entropy to the real sequence-level distribution. Another intuition behind the cross-entropy loss is that we want to find model parameters Θ that make the model distribution $P_{\Theta}(\cdot|\mathbf{x})$ similar to the *real* distribution $P(\cdot|\mathbf{x})$ over translations for a source sentence \mathbf{x} . The similarity is measured with the cross-entropy $H_{\mathbf{x}}(P, P_{\Theta})$. In practice, the real distribution $P(\cdot|\mathbf{x})$ is not known, but we have access to a training corpus of pairs (\mathbf{x}, \mathbf{y}) . For each such pair we consider the target sentence \mathbf{y} as a *sample* from the real distribution $P(\cdot|\mathbf{x})$. We now approximate the cross-entropy $H_{\mathbf{x}}(P, P_{\theta})$ using Monte Carlo estimation with only one sample ($N = 1$):

$$\begin{aligned} H_{\mathbf{x}}(P, P_{\Theta}) &= \mathbb{E}_{\mathbf{y}}[-\log P_{\Theta}(\mathbf{y}|\mathbf{x})] \\ &= - \sum_{\mathbf{y}'} P(\mathbf{y}'|\mathbf{x}) \log P_{\Theta}(\mathbf{y}'|\mathbf{x}) \\ &\stackrel{\text{MC}}{\approx} - \frac{1}{N} \sum_{\mathbf{y}'} \log P_{\Theta}(\mathbf{y}'|\mathbf{x}) \\ &\stackrel{N=1}{=} - \log P_{\Theta}(\mathbf{y}|\mathbf{x}) \\ &= - \sum_{j=1}^{|\mathbf{y}|} \log P_{\Theta}(y_j | y_1^{j-1}, \mathbf{x}) \\ &= \mathcal{L}_{\text{CE}}(\mathbf{x}, \mathbf{y}, \Theta). \end{aligned}$$

¹⁴I thank Adri   de Gispert for making that point in our discussions.

Cross-entropy loss optimizes a Monte Carlo approximation of the cross-entropy to the real token-level distribution. We arrive at the same result if we consider the cross-entropy between the *conditionals* of $P(\cdot|y_1^{j-1}, \mathbf{x})$ and $P_\Theta(\cdot|y_1^{j-1}, \mathbf{x})$ for given \mathbf{x} and prefix y_1^{j-1} :

$$\begin{aligned} \mathbb{E}_{y_j}[-\log P_\Theta(y_j|y_1^{j-1}, \mathbf{x})] &= -\sum_{y'_j} P(y'_j|y_1^{j-1}, \mathbf{x}) \log P_\Theta(y'_j|y_1^{j-1}, \mathbf{x}) \\ &\stackrel{\text{MC with } N=1}{\approx} -\sum_{j=1}^{|\mathbf{y}|} \log P_\Theta(y_j|y_1^{j-1}, \mathbf{x}) \\ &= \mathcal{L}_{\text{CE}}(\mathbf{x}, \mathbf{y}, \Theta). \end{aligned}$$

Cross-entropy loss optimizes the cross-entropy to the Dirac distribution. Alternatively, we can define a (Dirac) distribution which assigns the probability of one to \mathbf{y} and zero to all other target sentences:

$$P_\delta(\mathbf{y}'|\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{y}' = \mathbf{y} \\ 0 & \text{if } \mathbf{y}' \neq \mathbf{y} \end{cases} \quad (38)$$

The cross-entropy between the Dirac distribution (in this context taking the role of the empirical distribution) and our model distribution $P_\Theta(\cdot|\mathbf{x})$ is:

$$H_{\mathbf{x}}(P_\delta, P_\Theta) = -\sum_{\mathbf{y}'} P_\delta(\mathbf{y}'|\mathbf{x}) \log P_\Theta(\mathbf{y}'|\mathbf{x}) = -\log P_\Theta(\mathbf{y}|\mathbf{x}) \stackrel{\text{Eq. 37}}{=} \mathcal{L}_{\text{CE}}(\mathbf{x}, \mathbf{y}, \Theta). \quad (39)$$

To recap, we have found that the following are equivalent:

- Training under cross-entropy loss (Eq. 36).
- Maximizing the likelihood of the training data.
- Minimizing an estimate of the cross-entropy to the real sequence-level distribution.
- Minimizing an estimate of the cross-entropy to the real token-level distribution.
- Minimizing the cross-entropy to the Dirac distribution.

In particular, we emphasize the equivalence between the sequence-level and the token-level estimation since cross-entropy loss is often characterized as token-level objective in the literature whereas the term *sequence-level training* somewhat misleadingly usually refers to risk-based training under BLEU [318, 319] which is discussed in Sec. 11.5.

11.2. Training Deep Architectures

Deep encoders and decoders consisting of multiple layers have now superseded earlier shallow architectures. However, since the gradients have to be back-propagated through more layers, deep architectures – especially recurrent ones – are prone to vanishing gradients [320] and are thus harder to train. A number of tricks have been proposed recently that make it possible to train deep NMT models reliably. Residual connections [130] are direct connections that bypass more complex sub-networks in the layer stack. For example, all the architectures presented in Sec. 6.6 (GNMT, ConvS2S, Transformer,

RNMT+) add residual connections around attentional, recurrent, or convolutional cells to ease learning (Fig. 13). Another technique to counter vanishing gradients is called *batch normalization* [128] which normalizes the hidden activations in each layer in a mini-batch to a mean of zero and a variance of 1. An extension of batch normalization which is independent of the batch size and is especially suitable for recurrent networks is called *layer normalization* [129]. Layer normalization is popular for training deep NLP models like the Transformer [76].

11.3. Regularization

Modern NMT architectures are vastly over-parameterized [151] to help training [321]. For example, a subword-unit-level Transformer in a standard “big” configuration can easily have 200-300 million parameters [126]. The large number of parameters potentially makes the model prone to *over-fitting*: The model fits the training data perfectly, but the performance on held-out data suffers as the large number of parameters allows the optimizer to marginally improve training loss at the cost of generalization as training proceeds. Techniques that aim to prevent over-fitting in over-parameterized neural networks are called *regularizers*. Perhaps the two simplest regularization techniques are L1 and L2 regularization. The idea is to add terms to the loss function that penalize the magnitude of weights in the network. Intuitively, such penalties draw many parameters towards zero and limit their significance. Thus, L1 and L2 effectively serve as soft constraint on the model capacity.

The three most popular regularization techniques for NMT are *early stopping*, *dropout*, and *label smoothing*. Early stopping can be seen as regularization in time as it stops training as soon as the performance on the development set does not improve anymore. Dropout [127] is arguably one of the key techniques that have made deep learning practical. Dropout randomly sets the activities of hidden and visible units to zero during training. Thus, it can be seen as a strong regularizer for simultaneously training a large collection of networks with extensive weight sharing.

Label smoothing has been derived for expectation-maximization training by Byrne [322], and has been applied to large-scale computer vision by Szegedy et al. [299]. Label smoothing changes the training objective such that the model produces smoother distributions. We have already established in Sec. 11.1 that standard cross-entropy training measures the distance of the output distribution to the Dirac distribution around the training sample. Label smoothing discounts the likelihood of the training sample and distributes some of the free probability mass among other hypotheses. In NMT, label smoothing is applied as cross-entropy loss to a smoothed distribution $Q(\cdot)$ on the token level:

$$\mathcal{L}_Q(\mathbf{x}, \mathbf{y}, \Theta) = - \sum_{j=1}^{|\mathbf{y}|} \sum_{y' \in \Sigma_{trg}} Q(y'|j, \mathbf{y}) \log P_{\Theta}(y'|y_1^{j-1}, \mathbf{x}). \quad (40)$$

The distribution $Q(\cdot)$ can take language modelling scores into account [323], but usually it is just a smoothed version of the Dirac distribution for the reference label:

$$Q_{\alpha}(y'|j, \mathbf{y}) = \begin{cases} \alpha & \text{if } y' = \mathbf{y}_j \\ \frac{1-\alpha}{|\Sigma_{trg}|-1} & \text{if } y' \neq \mathbf{y}_j \end{cases} \quad (41)$$

for some smoothing factor $\alpha \in (0, 1]$. Setting $\alpha = 1$ recovers the normal cross-entropy loss from Sec. 11.1.

While label smoothing makes intuitive sense for computer vision, applying it to neural sequence prediction in this way has objectionable side effects on the sequence level. Considering the probabilities $Q(\cdot)$ assigns to *full sequences*, we first note that $Q(\cdot)$ does *not* uniformly distribute the remaining probability mass among all other sequences. In fact, distributing it uniformly would result in infinitely small probabilities as there are infinitely many possible sequences. Interestingly, $Q(\cdot)$ does also not assign a fixed probability of α to the correct sequence \mathbf{y} :

$$Q_\alpha(\mathbf{y}) = \prod_{j=1}^{|\mathbf{y}|} Q_\alpha(y_j | y_1^{j-1}) = \prod_{j=1}^{|\mathbf{y}|} \alpha = \alpha^{|\mathbf{y}|}. \quad (42)$$

Since α is less than one, $Q_\alpha(\cdot)$ is sharper if the correct sequence \mathbf{y} is short, and smoother if it is long. Alternative loss functions that encourage smooth output distributions include explicit entropy penalization [324] and knowledge distillation (Sec. 16). A regularization effect can also be achieved by making the training data harder to fit by adding noise, e.g. via subword regularization [259], SwitchOut [283], or noisy back-translation [282] (see Secs. 8.3 and 9).

11.4. Large Batch Training

Another practical trick which is becoming increasingly feasible with the availability of multi-GPU training and large GPU memories is to use very large batch sizes. Large batch training can yield almost linear speed-ups [325] as the computation can be distributed across multiple GPUs. Even more importantly, gradients estimated on large batches are naturally less noisy than gradients from small batches, and can yield better overall convergence [76, 126, 161]. For example, distributing Transformer training across 16 (effective) GPUs can improve over single GPU training by two full BLEU points [126]. Smith et al. [326] argued that increasing the batch size during training can have a similar effect as learning rate decay. For a thorough and insightful discussion of large batch training we refer the reader to [325].

Previous studies [327, 328] on batch size were limited by the hardware since – in vanilla SGD – the training batch has to fit into the GPU memories. Saunders et al. [329] presented a technique called *delayed SGD* which sidesteps these limitations by decoupling the batch size limit from the available hardware.

11.5. Reinforcement Learning

Ranzato et al. [318] pointed out two weaknesses of standard MLE training in neural sequence models. First, there is a discrepancy between NMT training and decoding. During training, the correct target label y_{j-1} is used in the j -th time step. Obviously, during decoding, the correct labels are not available, so the previous (potentially wrong) output is fed back to the model. This is called ‘exposure bias’ [318] as the model is never exposed to its own mistakes during training. The exposure bias can be tackled by feeding back the ground-truth labels only at early training stages, but gradually switching to feeding back the previously produced target tokens instead as training progresses [330].

The second issue in NMT training pointed out by Ranzato et al. [318] is the mismatch between training loss function and evaluation metric. Training uses cross-entropy loss on the word-level, whereas the final evaluation metric is usually BLEU [296] which is defined on sentence- or document-level. Both of these problems can be tackled with reinforcement learning [318, 331]. In the standard terminology of reinforcement learning, an *agent* interacts with an *environment* via *actions*. A *policy* determines the action to pick depending on the environment. The goal is to learn a policy which maximizes the expected *reward*. In NMT, the agent is the NMT model that interacts with the environment consisting of the source sentence \mathbf{x} and the translation history y_1^{j-1} by picking actions (words) according to the policy $P(y_j|y_1^{j-1}, \mathbf{x})$.

The advantage of casting NMT as reinforcement learning problem is that the reward does not need to be differentiable, and thus can be any quality measure such as BLEU or GLEU [14]. However, training is computationally very expensive as it requires sampling or decoding during training [332]. Therefore, reinforcement learning is usually used to refine a model trained with cross-entropy [14]. However, even though reinforcement learning has yielded some gains in the past in isolated experiments, it is difficult to improve over stronger baselines with recent NMT architectures and back-translation [333]. Wu et al. [14] reported that their gains in BLEU from reinforcement learning were not reflected in the human evaluation. Other possible applications for reinforcement learning in neural sequence prediction include architecture search [334], adequacy-oriented learning [308], and simultaneous translation (Sec. 7.8). An alternative way to incorporate the BLEU metric into NMT training is via a minimum risk formulation [319, 335].

11.6. Dual Supervised Learning

Recall that NMT networks are trained to model the distribution $P(\mathbf{y}|\mathbf{x})$ over translations \mathbf{y} for a given source sentence \mathbf{x} . This training objective takes only one translation direction into account – from the source language to the target language. However, the chain rule gives us the following relation:

$$P(\mathbf{y}|\mathbf{x})P(\mathbf{x}) = P(\mathbf{x}, \mathbf{y}) = P(\mathbf{x}|\mathbf{y})P(\mathbf{y}). \quad (43)$$

Eq. 43 is often not satisfied when the two translation models $P(\mathbf{y}|\mathbf{x})$ and $P(\mathbf{x}|\mathbf{y})$ are trained independently. The dual supervised learning loss \mathcal{L}_{DSL} aims to correlate both translation directions as follows [289, 336]:

$$\mathcal{L}_{\text{DSL}} = (\log P(\mathbf{x}) + \log P(\mathbf{y}|\mathbf{x}) - \log P(\mathbf{y}) - \log P(\mathbf{x}|\mathbf{y}))^2. \quad (44)$$

An alternative way to incorporate both translation directions is the agreement-based approach of Cheng et al. [337].

11.7. Adversarial Training

Generative adversarial networks [338, GANs] have recently become extremely popular in computer vision. GANs were originally proposed as framework for training generative models. For example, in computer vision, a generative model G would generate images that are similar to the ones in the training corpus. The input to a classic GAN is noise which is sampled from a noise prior. The key idea of adversarial training is that G is trained to fool a discriminative model D . The discriminator D takes an image as input

and outputs the probability of the image coming from the real training corpus as opposed to being generated by G . G and D are jointly trained with opposing objectives: G tries to drive up the probability of D making a mistake whereas D aims to discriminate between real and fake images generated by G . GANs are particularly useful when they condition on some input (conditional GANs). For example, a GAN which conditions on a textual description of an image is able to synthesize an image for an unseen description at test time.

In computer vision, it is possible to back-propagate gradients through the synthetic image and thus train G and D jointly without approximations. The main challenge for applying GANs to text is that this is no longer possible since text consists of a variable number of discrete symbols. Therefore, most work on adversarial training in NLP relies on reinforcement learning to generate synthetic text samples [339–343] or directly operates on the hidden activations in G [344]. Besides some exploratory efforts [339–341], adversarial training for NLP and particularly NMT is still in its infancy and rather brittle [340, 345–347].

12. Explainable Neural Machine Translation

12.1. *Post-hoc Interpretability*

Explaining the predictions of deep neural models is hard because they consist of tens of thousands of neurons and millions of parameters. Therefore, explainable and interpretable deep learning is still an open research question [348–352]. *Post-hoc interpretability* refers to the idea of sidestepping the model complexity by treating it as a black-box and not trying to understand the inner workings of the model. Montavon et al. [351] defines post-hoc interpretability as follows: “A trained model is given and our goal is to understand what the model predicts (e.g. categories) in terms what is readily interpretable (e.g. the input variables)”. In NMT, this means that we try to understand the target tokens (“what the model predicts”) in terms of the source tokens (“the input variables”). Post-hoc interpretability methods such as layer-wise relevance propagation [353] are often visualized with heat maps representing the importance of input variables – pixels in computer vision or source words in machine translation.

Applying post-hoc interpretability methods to sequence-to-sequence prediction has received some attention in the literature [354]. Alvarez-Melis and Jaakkola [355] proposed a causal model which finds related source-target pairs by feeding in perturbed versions of the source sentence. Ma et al. [356] derived relevance scores for NMT by comparing the predictive probability distributions before and after zeroing out a particular source word. See [357] for some general limitations of such post-hoc analyses in NLP.

12.2. *Model-intrinsic Interpretability*

Unlike the black-box methods for post-hoc interpretability, another line of research tries to understand the functions of individual hidden neurons or layers in the NMT network. Different methods have been proposed to visualize the activities or gradients in hidden layers [358–361]. Belinkov et al. [362] shed some light on NMT’s ability to handle morphology by investigating how well a classifier can predict part-of-speech or morphological tags from the last encoder hidden layer. Bau et al. [363], Dalvi et al. [364, 365] found individual neurons that capture certain linguistic properties with different

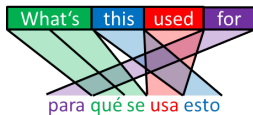


Figure 19: Word alignment from the English sentence “What’s this used for” to the Spanish sentence “para que se usa esto”.

forms of regression analysis. Bau et al. [363] were even able to alter the translation (e.g. change the gender) by manipulating the activities in these neurons. Other researchers have focused on the attention layer. Tang et al. [366] suggested that attention at different layers of the Transformer serves different purposes. They also showed that NMT does not use the means of attention for word sense disambiguation. Ghader and Monz [367] provide a detailed analysis of how NMT uses attention to condition on the source sentence.

12.3. Confidence Estimation in Translation

Obtaining word level or sentence level confidence scores for translations is not only very useful for practical MT, it also improves the explainability and trustworthiness of the MT system. An obvious candidate for confidence scores from an NMT system are the probabilities the model assigns to tokens or sentences. However, there is some disagreement in the literature on how well NMT models are calibrated [298, 368]. Poorly calibrated models do not assign probabilities according to the true data distribution. Such models might still assign high scores to high quality translations, but their output distributions are no reliable source for deriving word-level confidence scores. While confidence estimation has been explored for traditional SMT [369–371], it has received almost no attention since the advent of neural machine translation. The only work on confidence in NMT we are aware of is from Rikters and Fishel [372], Rikters [373] who aim to use attention to estimate word-level confidences.

In contrast, the related field of Quality Estimation for MT enjoys great popularity, with well-attended annual WMT evaluation campaigns – by now in their seventh edition [374]. Quality estimation aims to find meaningful quality metrics which are more accepted by users and customers than abstract metrics like BLEU [296], and are more correlated to the usefulness of MT in a real-world scenario. Possible applications for quality estimation include estimating post-editing efficiency [375] or selecting sentences in the MT output which need human revision [370].

12.4. Word Alignment in Neural Machine Translation

Word alignment is one of the fundamental problems in traditional phrase-based SMT. SMT constructs the target sentence by matching phrases in the source sentence, and combining their translations to form a fluent sentence [181, 217]. This approach does not only yield a translation, it also produces a word alignment along with it since each target phrase is generated from a unique source phrase. Thus, a word alignment can be seen as an explanation for the produced translation: each target phrase is explained with a link into the source sentence (Fig. 19). Unfortunately, vanilla NMT does not have the notion of a hard word alignment. It is tempting to interpret encoder-decoder attention matrices in neural models (Sec. 6.1) as (soft) alignments, but previous work has found that the

attention weights in NMT are often erratic and differ significantly from traditional word alignments:

- “The attention model for NMT does not always fulfill the role of a word alignment model, but may in fact dramatically diverge.” [300]
- “We perform extensive experiments across a variety of NLP tasks that aim to assess the degree to which attention weights provide meaningful ‘explanations’ for predictions. We find that they largely do not.” [376]
- “Attention weights are only noisy predictors of even intermediate components’ importance, and should not be treated as justification for a decision.” [377]
- “Although attention is very useful for understanding the connection between source and target words, only using attention is not sufficient for deep interpretation of target word generation.” [360]
- “Attention agrees with traditional alignments to a high degree in the case of nouns. However, it captures other information rather than only the translational equivalent in the case of verbs.” [367]
- “Attention visualizations are misleading and should be treated with care when explaining the underlying deep learning system.” [378]

Despite considerable consensus about the importance of word alignments for practical machine translation [300], e.g. to enforce constraints on the output [379] or to preserve text formatting, introducing explicit alignment information to NMT is still an open research problem. Word alignments have been used as supervision signal for the NMT attention model [380–383]. Cohn et al. [384] showed how to reintroduce concepts known from traditional statistical alignment models [270] like fertility and agreement over translation direction to NMT.

Hard attention [82] is a discrete version of the usual soft attention and is thus closer to the concept of a hard alignment. Similar ideas have been explored for speech recognition [385], morphological inflection [386], text summarization [387, 388], and image caption generation [82]. Some approaches to simultaneous translation presented in Sec. 7.8 explicitly control for reading source tokens and writing target tokens and thereby generate monotonic hard alignments on the segment level [210, 389]. Hybrids between soft and hard attention have been proposed by Choi et al. [390], Shen et al. [391]. However, the usefulness of hard attention for generic offline machine translation is often limited since it usually can only represent monotonic alignments.

Alkhouli et al. [392] used separate alignment and lexical models and thus were able to hypothesize explicit alignment links during decoding. Alignment-based NMT has been extended to multi-head attention by using an additional alignment head [393]. A similar idea was pursued by Zenkel et al. [394] who added an additional alignment layer to the Transformer and trained it – unlike Alkhouli et al. [393] – in an unsupervised way. The neural operation sequence model of Stahlberg et al. [171] is another way of generating an alignment along with the translation in NMT.

13. Alternative NMT Architectures

13.1. Extensions to the Transformer Architecture

The Transformer model architecture [76] introduced in Sec. 6.5 has become the de facto standard architecture for neural machine translation because of its superior translation quality on a variety of language pairs [145, 146].¹⁵ The Transformer comes with a number of techniques which sets it apart from previous architectures such as multi-head attention, self-attention, large batch training, etc. Some ablation studies in the literature aim to factor out or explain the contributions of these different techniques [123, 124, 366, 396]. Several attempts have been made to improve different aspects of the vanilla model for machine translation, but none has been widely adopted. Most notably, Shaw et al. [125] proposed to embed relative positions rather than absolute ones. A disadvantage of the relative Transformer is the increased computational complexity. The memory keys and values with absolute positions are the same in each decoding step. With relative positioning, however, both have to be recomputed in each time step since the relative positions change over time. The model of Song et al. [397] works with attention masks (Sec. 6.2) to narrow down context. Ahmed et al. [398] proposed to weight the output of attention heads inside multi-head attention. The Star-Transformer [399] thins out inter-layer connections of the standard model to reduce computational complexity. With a similar outset, Medina and Kalita [400] reported speed-ups by replacing the single deep encoder with multiple shallow encoders.

Some recent research has focused on large scale language modelling with the Transformer [48, 132, 401–403]. The Transformer is also the starting point for neural architectures for contextualized word embeddings (see Sec. 2) such as BERT [49].

13.2. Advanced Attention Models

As shown in Sec. 6.1, the vast majority of current NMT architectures are based on one of three attention types: additive, (scaled) dot-product, or multi-head attention [13, 76, 77]. In this section, we will outline attempts to improve upon these standard models.

Sec. 10.1 discussed the problem of over- and under-translation, and how coverage models can mitigate this problem by controlling the attention weights with fertilities. Alternatively, researchers have tried to equip the attention layer itself with additional components like a memory [404] or a recurrent network [405, 406] to enable it to keep track of the attention history. Choi et al. [407] proposed an attention model that is able to learn different attention weights for each dimension in the values, not only one weight for each value vector.

One potential weakness of the standard models is that they are token-based: the attention output is a weighted average of the values, and the attention weights tend to focus on a single key-value pair. Therefore, there is no explicit mechanism to attend to full phrases rather than subwords or characters.¹⁶ Phrase-based NMT which equips the model with the ability to attend to full phrases or multi-word expressions has been

¹⁵The only contrary evidence we are aware of is from Tran et al. [395] who found that recurrent models can better model hierarchical structure than the Transformer.

¹⁶This does not mean that the source sentence context is always reduced to a single input token since the encoder hidden states are by themselves context-sensitive.

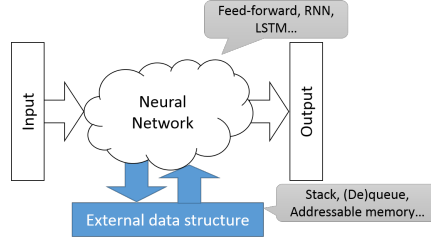


Figure 20: Neural networks with external memory.

studied by Rikters and Bojar [408], Ishiwatari et al. [409], Feng et al. [410], Huang et al. [411], Li et al. [412], Eriguchi et al. [413].

On the other side of the spectrum, it has been noted that regular attention sometimes spreads out over too many elements, especially when applied over long sequences. The attention output in this case is an average of many values which is naturally more noisy than with sharp attention, and which impedes the propagation of information through the network. Hard attention (Sec. 12.4) removes this sort of noise, but is often restricted to monotonic alignment. Lin et al. [414] proposed to explicitly learn to set the temperature of attention weights to control the softness of attention. Another potential solution has been suggested by Zhang et al. [415] who used GRU gates rather than weighted linear combinations to compute the attention output from the values.

13.3. Memory-augmented Neural Networks

RNNs are theoretically Turing-complete [416] and thus potentially very powerful models of computation. However, since training is still a challenge (see Sec. 11), even advanced RNN architectures like LSTMs [71] fail to solve certain basic sequence-to-sequence tasks like (repeated) copying or reversal in practice [417, 418]. This observation motivated researchers to add external memory structures like a memory tape [419] or a stack [420, 421] to the neural network. The basic idea is illustrated in Fig. 20. Besides producing the output sequence, the neural network learns to operate an external data structure. The external memory is not part of the neural network but the network learns to communicate with it through conceptually discrete operations like PUSH and POP. However, in order to train the whole system with a gradient-based optimizer, these discrete operations are often approximated with continuous versions [417, 418, 422]. Various data structures have been used in combination with neural networks such as (inter alia) stacks [422], (double-ended) queues [418], addressable memory cells [417, 423, 424], and hierarchical memory structures [425]. Grefenstette et al. [418] suggested that even simple data structures like dequeues help to solve linguistically motivated tasks like bigram flipping or Inversion Transduction Grammar [55, ITG] tasks. Research on these kinds of neural network operated data structures still mainly focuses on synthetic tasks like relatively simple algorithmic problems. Initial efforts to apply this line of research to real world problems are limited to neural machine translation [426–429], sentence simplification [430], and text normalization [431].

13.4. Beyond Encoder-decoder Networks

All NMT architectures which we have discussed in the previous sections fall in the category of encoder-decoder networks: An encoder network computes a fixed or variable length continuous hidden representation of the source sentence, and a separate decoder network defines a probability distribution over target sentences given that representation. There are some initial efforts in the literature to depart from this overall structure. For example, variational methods that define a *distribution* over (a part of) the hidden representations have been explored by Zhang et al. [432], Su et al. [433], Bastings et al. [434], Shah and Barber [435]. Non-autoregressive NMT which aims to reduce or remove the sequential dependency on the translation prefix inside the decoder for enhanced parallelizability has been studied by Wang et al. [436], Gu et al. [437], Guo et al. [438], Wang et al. [439], Libovický and Helcl [440], Lee et al. [441], Akoury et al. [442]. Bahar et al. [443], Kaiser and Bengio [444] recomputed the encoder state after each time step and thus effectively expanded the hidden representation into a 2D structure. The architecture proposed by He et al. [445] does not only use the last encoder layer as hidden representation, but instead connects encoder and decoder layers at the same depth via attention.

14. Data Sparsity

Deep learning methods are notoriously data hungry. For example, traditional statistical machine translation still often outperforms neural machine translation when training data is scarce [169, 300]. In this section we will look at the problem of training data sparsity from different angles such as reducing noise in training data (Sec. 14.1), using data from a different domain, or making use of less or no parallel data.

14.1. Corpus Filtering

Unfortunately, MT training data is usually inherently noisy as it is often extracted (semi-) automatically by crawling the web [446, 447] and therefore commonly contains sentence fragments, wrong languages, misaligned sentence pairs [448], or MT output rather than genuine parallel text [449, 450]. In the previous sections we discussed several instances of the use of synthetic noise in NMT. For example, adding noise to the synthetic sentences in back-translation can be beneficial (Sec. 9). Noise can also be used to generate diverse translations (Sec. 7.7) or as regularizer (Sec. 11.3). However, when discussing the role of noise in NMT it is imperative to carefully differentiate between the various kinds of noise and the ways it impacts NMT. Studies have shown that NMT is not robust against naturally occurring noise at training [448] and test [268, 451–453] time. Robustness at test time can be improved by training on synthetic noise [454, 455]. Corpus filtering to reduce the amount of noise in the training data has been widely studied for traditional SMT [456, 457], often in context of domain adaptation [458, 459]. More recent research on data filtering focuses on NMT since van der Wees et al. [460] had shown that filtering techniques developed for SMT are less useful for NMT. One of the first approaches to NMT corpus filtering was the method of Carpuat et al. [461] based on semantic analysis. The most effective approaches in the WMT18 shared task on corpus filtering for NMT [462] used a combination of likelihood scores from neural translation models and neural language models which have been trained on clean data [149, 463,

464]. These criteria prefer sentence pairs which are likely translations of one another according to the translation model [465]. Zhang et al. [466] proposed the exact opposite, arguing that NMT training should concentrate on “difficult” training samples, i.e. samples with low translation probability. An alternative to hard data filtering called *curriculum learning* [467] that controls the order of training samples has been applied to NMT by van der Wees et al. [460], Wang et al. [468], Kumar et al. [469], Platanios et al. [470].

14.2. Domain Adaptation

There is a robust body of research on domain adaptation for machine translation [471, 472]. Popular domain adaptation techniques for both SMT and NMT aim to select [458, 459, 473–475] or weight [474, 476, 477] samples in a large out-of-domain corpus. Back-translation (Sec. 9) can also be used for domain adaptation by back-translating sentences from an in-domain monolingual corpus. Another simple yet very effective method is to jointly train on in-domain and out-domain sentences, possibly with domain-tags to help learning [478–480]. Sajjad et al. [481] showed that a simple concatenation of in-domain and out-domain corpora can already increase the robustness and generalization of NMT significantly. Khayrallah et al. [482] studied domain adaptation by constraining an NMT system to SMT lattices. Freitag and Al-Onaizan [152] ensembled separately trained general-domain and in-domain models.

Another widely used technique is to train the model on a general domain corpus, and then fine-tune it by continuing training on the in-domain corpus [274, 483]. Fine-tuning bears the risk of two negative effects: catastrophic forgetting [484, 485] and over-fitting. Catastrophic forgetting occurs when the performance on the specific domain is improved after fine-tuning, but the performance of the model on the general domain has decreased drastically. The risk of over-fitting is connected to the fact that the in-domain corpus is usually very small. Both effects can be mitigated by artificially limiting the learning capabilities of the fine-tuning stage, e.g. by freezing sub-networks [486] or by only learning additional scaling factors for hidden units rather than full weights [487, 488]. A very elegant way to prevent over-fitting and catastrophic forgetting is to apply regularizers (Sec. 11.3) to keep the adapted model weights close to their original values. Khayrallah et al. [489], Dakwale and Monz [490] regularized the output distributions using techniques inspired by knowledge distillation (Sec. 16). Miceli Barone et al. [491] applied standard L2 regularization and a variant of dropout to domain adaptation. Elastic weight consolidation [492] can be seen as generalization of L2 regularization that takes the importance of weights (in terms of Fisher information) into account, and has been applied to NMT domain adaptation by Thompson et al. [493], Saunders et al. [494]. In particular, Saunders et al. [494] showed that EWC does not only reduce catastrophic forgetting but even yields gains on the general domain when used for fine-tuning on a related domain.

14.3. Low-resource NMT

One of the areas in which traditional SMT still often outperforms NMT is low-resource translation [169, 300]. However, several techniques have been proposed to improve the performance of NMT under low-resource conditions. In general, the methods discussed in Sec. 9 to leverage monolingual data such as back-translation are particularly effective for low-resource MT. Ren et al. [495] proposed a scheme that could make use of translations from/into the source/target language into/from a third resource-rich language. The

transfer-learning approach of Zoph et al. [496] first trains a *parent* model on a resource-rich language pair (e.g. French-English), and then continues training on the low-resource pair of interest (e.g. Uzbek-English). The effectiveness of transfer-learning depends on the relatedness of the languages [496–499]. The rapid adaptation of multilingual NMT systems to new low-resource language pairs has been studied by Neubig and Hu [500]. Approaches that do not rely on resources from a third language include Östling and Tiedemann [169] who supervised the generation order of an insertion-based low-resource translation model with word alignments.

A series of NIST evaluation campaigns called LoReHLT [501] focuses on low-resource MT, and recent WMT editions also contain low-resource language pairs [144–146].

14.4. Unsupervised NMT

Unsupervised NMT is an extreme case of the low-resource scenario in which not even small amounts of cross-lingual data is available, and the translation system learns entirely from (unrelated) monolingual data. Unsupervised NMT often starts off from an unsupervised cross-lingual word embedding model [502–504] that maps word embeddings from the source and the target language into a joint embedding space [505, 506]. The translation model is then further refined by iterative back-translation [507, 508]. The extract-edit scheme of Wu et al. [509] is an alternative to back-translation for unsupervised NMT that edits a sentence in the monolingual corpus rather than synthesize it from scratch. Unsupervised NMT has been targeted in recent WMT evaluation campaigns [145, 146].

15. Multilingual NMT

NMT is usually trained to translate a single fixed source language into another fixed target language. Multilingual NMT aims to cover translation directions between multiple languages with a single model. This does not only have the potential of exploiting similarities across language pairs, it also reduces the number of systems required for all-way translation between a set of languages from quadratic to linear or even one. Multilingual NMT systems can be largely categorized by the components they share between language directions. On one side of the spectrum, the entire neural architecture (both encoder and decoder) can be shared, and source and target languages can be specified by annotating sentences [510] or words [511, 512] with language ID tags or embeddings. On the other side of the spectrum, Luong et al. [513] used a separate encoder for each source language and a separate decoder for each target language. Firat et al. [514, 515] extended the work of Luong et al. [513] to attentional NMT by sharing the attention mechanism across language directions. Dong et al. [516] studied one-to-many translation with a single encoder but separate decoders for each target language. A potential benefit of multilingual systems is zero-shot translation, i.e. the translation between two languages for which no direct training data is available.¹⁷ Johnson et al. [510] reported reasonable Portuguese→Spanish translation performance of their multilingual system

¹⁷The difference between zero-shot and unsupervised NMT (Sec. 14.4) is that unsupervised NMT does not rely on *any* cross-lingual data whereas zero-shot NMT uses cross-lingual data in other language directions.

that has been trained on Portuguese \leftrightarrow English and Spanish \leftrightarrow English, although pivoting through English (translate Spanish to English, and then English to Portuguese) worked better. Pivot-based zero-shot translation can be further improved by fine-tuning on a pseudo parallel corpus [154] or by jointly training some components of the source-pivot and pivot-target systems like word embedding matrices [517]. Lu et al. [518] reported gains in zero-shot settings by adding a boldly named “neural interlingual” component between the encoder and the decoder which is shared across language directions. For an assessment of the current capabilities of multilingual and zero-shot translation systems see [519–521]. Another form of multilingual NMT is multi-source NMT [155, 522], in which the system tries to generate a single translation given sentences in two source languages simultaneously. A problem with this approach is data sparsity as missing source sentences have to be synthesized [523, 524] if the training corpus does not provide sentences in all source languages. In a wider context, multi-source architectures can be used for multimodal NMT (Sec. 17.1), morphological inflection [525], zero-shot translation [154], low-resource MT [523], syntax-based NMT [526], document-level MT [527], or bidirectional decoding [172]. Dabre et al. [528] provide an overview of recent trends in multilingual NMT.

16. NMT Model Size

NMT models usually have hundreds of millions of parameters (Tab. 4). Such large models cause a number of practical issues. GPUs are usually required to run such big models efficiently, but GPUs are expensive and their memory is limited. Smaller models would not only reduce the computational complexity but could also make better use of GPU parallelism by increasing batch sizes. Furthermore, model files require large amounts of disk space which is a problem on mobile platforms. One way to increase the space efficiency of neural models is neural architecture search [132, 334]. For example, So et al. [529] found computationally efficient Transformer hyper-parameters by systematic neural architecture search. Rather than optimizing the dimensionality of layers, it is also possible to significantly speed up translation by departing from the usual 32 bit floating point arithmetics by reducing the precision to 8 or 16 bits [530–533] or by using vector quantization [14, 534]. The idea of pruning neural networks to improve the compactness of the models dates back almost 30 years [535]. The literature is therefore vast [536]. One line of research aims to remove unimportant network connections. The connections can be selected for deletion based on the second-derivative of the training error with respect to the weight [535, 537], or by a threshold criterion on its magnitude [538]. See et al. [539] confirmed a high degree of weight redundancy in NMT networks. Zhu and Gupta [540] demonstrated that large sparse models outperform smaller dense networks with the same memory footprint. Srinivas and Babu [541] proposed to remove neurons which are very similar to another neuron and have small outgoing weights. Stahlberg and Byrne [151] generalized their method to linear combinations of neurons. Babaeizadeh et al. [542] combined pairs of neurons with similar activities during training. Using low rank matrices for neural network compression, particularly approximations via Singular Value Decomposition (SVD), has been studied widely in the literature [543–547]. Another approach, known as *knowledge distillation*, uses a large model (the teacher) to generate soft training labels for a smaller student network [548, 549]. The student network is trained by minimizing the cross-entropy to the teacher. This idea has been applied to

sequence modelling tasks such as machine translation and speech recognition [190, 550–554].

17. NMT with Extended Context

17.1. Multimodal NMT

Machine translation is usually framed as the isolated transformation of the textual representation of a single sentence in one language into another. Since language is inherently ambiguous, researchers have searched for ways to provide the translation system with more context. For example, if the source sentence describes an image, the image itself potentially carries valuable clues to help the translation process. Multimodal machine translation [555, 556] aims to generate an image caption in the target language given both the source language caption and the image itself. The core of most multimodal MT models is a normal text-to-text system which integrates visual information by using global image features extracted with a separate computer vision model [555, 556] or via visual attention [557]. Multimodality in translation was the subject of a series of WMT shared tasks [558–560]. Calixto and Liu [561] demonstrated the usefulness of visual clues in translation.

17.2. Tree-based NMT

The prevalent choice for modeling units in NMT are characters or subword-units (Sec. 8.3). This design decision is not linguistically motivated but rather stems from the difficulty of extending NMT to an open vocabulary. From the linguistic perspective, however, translation is better viewed as the transformation of larger elements in the sentence such as words, phrases, or even syntactic structures.

Various attempts have been made to introduce structures such as syntactic constituency trees or dependency trees both on the source and the target side of NMT. A popular approach is to retain the sequence-to-sequence architecture and linearize the tree structures, for example using bracket expressions [526, 562–564], sequences of rules [329], or CCG supertags [565]. Ma et al. [566], Zaremoondi and Haffari [567] developed a linearization of a packed forests that represented multiple source sentence parses. Saunders et al. [329] reported gains by ensembling different linearization strategies of target-side syntax trees. Recurrent neural network grammars [568] that represent syntactic parse trees as sequence of actions were applied to machine translation by Eriguchi et al. [569], Bradbury and Socher [570]. Using actions to build target side tree structures is also central to the tree-based decoders of Wang et al. [571], Wu et al. [572]. Akoury et al. [442] used syntax to speed up decoding by first predicting a parse tree, and then predicting all target tokens in parallel. Tree-LSTMs [573] make it possible to represent a tree structure directly with the neural network architecture. They are a generalization of recurrent LSTM cells (Sec. 6.3) that replaces the single input of a standard LSTM cell (usually from the previous time step) with multiple input connections, one from each child node. Thus, each Tree-LSTM cell represents a node in the tree, and the root node contains a fixed-length vector encoding of the whole tree structure. Tree-LSTMs have been applied to syntax-based NMT [574–576]. An alternative to Tree-LSTMs was proposed by Shen et al. [577] who rearranged neurons in an LSTM network to resemble

a block representation of the tree. Bastings et al. [578], Chen et al. [579] used convolutional encoders to represent a dependency graph in the source sentence. Chen et al. [580] biased encoder-decoder attention weights with syntactic clues. Unsupervised tree-based methods have been studied by Kim et al. [581], Maillard et al. [582], Williams et al. [583].

17.3. NMT with Graph Structured Input

As a generalization of the tree-based approaches discussed in the previous section, lattice-based NMT allows more general graph structures on the input side to provide a richer description of the source sentence. Lattices can represent uncertainty of upstream components such as speech recognizers [584] or tokenizers [585, 586]. Lattices have also been used to augment the input with external knowledge sources such as knowledge graphs [587, 588] or semantic predicate-argument structures [589]. Factors are another way of providing more information to the translation system. Factors describe a word by a tuple consisting of its lemma and various linguistic information (prefix, suffix, part-of-speech etc.) rather than its surface form. This technique is popular for traditional statistical machine translation [181, 590], and has been applied to neural machine translation both on the input [591] and the output [592, 593] side.

17.4. Document-level Translation

MT systems usually translate sentences in isolation. However, there is evidence that humans also take context into account, and rate translations from humans with access to the full document higher than the output of a state-of-the-art sentence-level machine translation system [594]. Common examples of ambiguity which can be resolved with cross-sentence context are pronoun prediction or coherency in lexical choice.

Various techniques have been proposed to provide the translation system with inter-sentential context, for example by initializing encoder or decoder states [595], using multi-source encoders [527, 596], as additional decoder input [595], with memory-augmented neural networks [597–599], a document-level LM [600], hierarchical attention [601, 602], deliberation networks [603], or by simply concatenating multiple source and/or target sentences [527, 604]. Context-aware extensions to Transformer encoders have been proposed by Voita et al. [605], Zhang et al. [606]. Techniques also differ in whether they use source context only [595, 596, 600, 605, 606], target context only [597, 599], or both [527, 598, 601, 602, 604]. Several studies on document-level NMT indicate that automatic and human sentence-level evaluation metrics often do not correlate well with improvements in discourse level phenomena [527, 594, 607].

18. NMT-SMT Hybrid Systems

Neural models were increasingly used as features in traditional SMT until NMT evolved as new paradigm. Without question, NMT has become the prevalent approach to machine translation in recent years. There is a large body of research comparing NMT and SMT (Tab. 6). Most studies have found superior overall translation quality of NMT models in most settings, but complementary strengths of both paradigms. Therefore, the literature about hybrid NMT-SMT systems is also vast. We distinguish between two categories of approaches for blending SMT and NMT.

Neural machine translation	Statistical machine translation
<ul style="list-style-type: none"> + Much better overall translation quality than SMT with enough training data [28, 300, 608–612]. + More fluent than SMT [608, 609, 611, 613, 614]. + Better handles a variety of linguistic phenomena than SMT [609, 610, 615]. – Adequacy issues due to lack of explicit coverage mechanism [306–308, 613, 614]. – Lack of hypothesis diversity (Sec. 7.7). – Neural models perform not as well as specialized symbolic models on several monotone seq2seq tasks [616]. 	<ul style="list-style-type: none"> + Outperforms NMT in low-resource scenarios [300, 613, 617–620]. + Produces richer output lattices [235]. + More robust against noise [448, 453]. + Translation quality degrades less on very long sentences than NMT [608, 609]. + Less errors in the translation of proper nouns [610]. <ul style="list-style-type: none"> ◦ NMT and SMT require comparable amounts of (document-level) post-editing [611, 621].

Table 6: Summary of studies comparing traditional statistical machine translation and neural machine translation.

Approaches in the first category do not employ a full SMT system but borrow only key ideas or components from SMT to address specific issues in NMT. It is straight-forward to combine NMT scores with other features normally used in SMT (like language models) in a log-linear model [271, 303].¹⁸ Conventional symbolic SMT-style lexical translation tables can be incorporated into the NMT decoder by using the soft alignment weights of the standard NMT attention model [139, 303, 622–624]. Cohn et al. [384] proposed to enhance the attention model in NMT by implementing basic concepts from the original word alignment models [270, 625] like fertility and relative distortion.

The second category of hybrid systems is related to system combination. The idea is to combine a fully trained SMT system with an independently trained NMT system. Popular examples in this category are rescoring and reranking methods [235, 482, 626–630], although these models may be too constraining if the neural system is much stronger. Stahlberg et al. [631] proposed a finite state transducer based loose combination scheme that combines NMT and SMT translations via an edit distance based loss. The minimum Bayes risk (MBR) based approach of Stahlberg et al. [199] biases an unconstrained NMT decoder towards n -grams which are likely according the SMT system, and therefore also does not constrain the system to the SMT search space. MBR-based combination of NMT and SMT has been used in WMT evaluation systems [126, 600] and in the industry [198]. NMT and SMT can also be combined in a cascade, with SMT providing the input to a post-processing NMT system [632, 633] or vice versa [634]. Wang et al. [635, 636] interpolated NMT posteriors with word recommendations from SMT and jointly trained NMT together with a gating function which assigns the weight between SMT and NMT scores dynamically. The AMU-UEDIN submission to WMT16 let SMT take the lead

¹⁸Note that this is still different from using neural features in an SMT system as the standard left-to-right NMT decoder is used.

and used NMT as a feature in phrase-based MT [159]. In contrast, Long et al. [637] translated most of the sentence with an NMT system, and just used SMT to translate technical terms in a post-processing step. Dahlmann et al. [638] proposed a hybrid search algorithm in which the neural decoder expands hypotheses with phrases from an SMT system. SMT can also be used as regularizer in unsupervised NMT [508].

19. Conclusion

Neural machine translation (NMT) has become the de facto standard for large-scale machine translation in a very short period of time. This article traced back the origin of NMT to word and sentence embeddings and neural language models. We reviewed the most commonly used building blocks of NMT architectures – recurrence, convolution, and attention – and discussed popular concrete architectures such as RNNsearch, GNMT, ConvS2S, and the Transformer. We discussed the advantages and disadvantages of several important design choices that have to be made to design a good NMT system with respect to decoding, training, and segmentation. We then explored advanced topics in NMT research such as explainability and data sparsity.

References

- [1] Y. Goldberg, A primer on neural network models for natural language processing, *Journal of Artificial Intelligence Research* 57 (2016) 345–420.
- [2] Y. Bengio, R. Ducharme, P. Vincent, C. Jauvin, A neural probabilistic language model, *Journal of machine learning research* 3 (2003) 1137–1155.
- [3] Y. Bengio, H. Schwenk, J.-S. Senécal, F. Morin, J.-L. Gauvain, *Neural Probabilistic Language Models*, Springer Berlin Heidelberg, Berlin, Heidelberg, 2006, pp. 137–186. URL: https://doi.org/10.1007/3-540-33486-6_6. doi:10.1007/3-540-33486-6_6.
- [4] H. Schwenk, D. Dechelotte, J.-L. Gauvain, Continuous space language models for statistical machine translation, in: *Proceedings of the COLING/ACL 2006 Main Conference Poster Sessions*, Association for Computational Linguistics, Sydney, Australia, 2006, pp. 723–730. URL: <https://www.aclweb.org/anthology/P06-2093>.
- [5] F. Zamora-Martinez, M. J. Castro-Bleda, H. Schwenk, N-gram-based machine translation enhanced with neural networks for the French-English BTEC-IWSLT’10 task, in: *International Workshop on Spoken Language Translation (IWSLT)* 2010, 2010.
- [6] H. Schwenk, Continuous space translation models for phrase-based statistical machine translation, in: *Proceedings of COLING 2012: Posters*, The COLING 2012 Organizing Committee, Mumbai, India, 2012, pp. 1071–1080. URL: <https://www.aclweb.org/anthology/C12-2104>.
- [7] H.-S. Le, A. Allauzen, F. Yvon, Continuous space translation models with neural networks, in: *Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Association for Computational Linguistics, Montréal, Canada, 2012, pp. 39–48. URL: <https://www.aclweb.org/anthology/N12-1005>.
- [8] J. Devlin, R. Zbib, Z. Huang, T. Lamar, R. Schwartz, J. Makhoul, Fast and robust neural network joint models for statistical machine translation, in: *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Association for Computational Linguistics, Baltimore, Maryland, 2014, pp. 1370–1380. URL: <https://www.aclweb.org/anthology/P14-1129>. doi:10.3115/v1/P14-1129.
- [9] N. Kalchbrenner, P. Blunsom, Recurrent continuous translation models, in: *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Seattle, Washington, USA, 2013, pp. 1700–1709. URL: <https://www.aclweb.org/anthology/D13-1176>.
- [10] K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, Y. Bengio, Learning phrase representations using RNN encoder-decoder for statistical machine translation, in: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*

- (EMNLP), Association for Computational Linguistics, Doha, Qatar, 2014, pp. 1724–1734. URL: <https://www.aclweb.org/anthology/D14-1179>. doi:10.3115/v1/D14-1179.
- [11] K. Cho, B. van Merriënboer, D. Bahdanau, Y. Bengio, On the properties of neural machine translation: Encoder–decoder approaches, in: Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation, Association for Computational Linguistics, Doha, Qatar, 2014, pp. 103–111. URL: <https://www.aclweb.org/anthology/W14-4012>. doi:10.3115/v1/W14-4012.
 - [12] I. Sutskever, O. Vinyals, Q. V. Le, Sequence to sequence learning with neural networks, in: Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, K. Q. Weinberger (Eds.), Advances in Neural Information Processing Systems 27, Curran Associates, Inc., 2014, pp. 3104–3112. URL: <http://papers.nips.cc/paper/5346-sequence-to-sequence-learning-with-neural-networks.pdf>.
 - [13] D. Bahdanau, K. Cho, Y. Bengio, Neural machine translation by jointly learning to align and translate, in: ICLR, 2015.
 - [14] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, et al., Google’s neural machine translation system: Bridging the gap between human and machine translation, arXiv preprint arXiv:1609.08144 (2016).
 - [15] J. Crego, J. Kim, G. Klein, A. Rebollo, K. Yang, J. Senellart, E. Akhanov, P. Brunelle, A. Coquard, Y. Deng, et al., SYSTRAN’s pure neural machine translation systems, arXiv preprint arXiv:1610.05540 (2016).
 - [16] T. Schmidt, L. Marg, How to move to neural machine translation for enterprise-scale programming—Tān early adoption case study, 2018.
 - [17] P. Levin, N. Dhanuka, T. Khalil, F. Kovalev, M. Khalilov, Toward a full-scale neural machine translation in production: the booking.com use case, arXiv preprint arXiv:1709.05820 (2017).
 - [18] G. Neubig, Neural machine translation and sequence-to-sequence models: A tutorial, arXiv preprint arXiv:1703.01619 (2017).
 - [19] F. Cromieres, T. Nakazawa, R. Dabre, Neural machine translation: Basics, practical aspects and recent trends, in: Proceedings of the IJCNLP 2017, Tutorial Abstracts, Asian Federation of Natural Language Processing, Taipei, Taiwan, 2017, pp. 11–13. URL: <https://www.aclweb.org/anthology/I17-5004>.
 - [20] P. Koehn, Neural machine translation, arXiv preprint arXiv:1709.07809 (2017).
 - [21] A. Popescu-Belis, Context in neural machine translation: A review of models and evaluations, arXiv preprint arXiv:1901.09115 (2019).
 - [22] A. Vaswani, S. Bengio, E. Brevdo, F. Chollet, A. N. Gomez, S. Gouws, L. Jones, Ł. Kaiser, N. Kalchbrenner, N. Parmar, R. Sepassi, N. Shazeer, J. Uszkoreit, Tensor2Tensor for neural machine translation, in: Proceedings of the 13th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Papers), Association for Machine Translation in the Americas, Boston, MA, 2018, pp. 193–199. URL: <https://www.aclweb.org/anthology/W18-1819>.
 - [23] M. Ott, S. Edunov, A. Baevski, A. Fan, S. Gross, N. Ng, D. Grangier, M. Auli, fairseq: A fast, extensible toolkit for sequence modeling, in: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations, 2019.
 - [24] G. Klein, Y. Kim, Y. Deng, J. Senellart, A. M. Rush, OpenNMT: Open-source toolkit for neural machine translation, in: Proceedings of ACL 2017, System Demonstrations, Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 67–72. URL: <https://www.aclweb.org/anthology/P17-4012>.
 - [25] F. Hieber, T. Domhan, M. Denkowski, D. Vilar, A. Sokolov, A. Clifton, M. Post, Sockeye: A toolkit for neural machine translation, arXiv preprint arXiv:1712.05690 (2017).
 - [26] O. Kuchaiev, B. Ginsburg, I. Gitman, V. Lavrukhin, C. Case, P. Micikevicius, OpenSeq2Seq: Extensible toolkit for distributed and mixed precision training of sequence-to-sequence models, in: Proceedings of Workshop for NLP Open Source Software (NLP-OSS), Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 41–46. URL: <https://www.aclweb.org/anthology/W18-2507>.
 - [27] R. Sennrich, O. Firat, K. Cho, A. Birch, B. Haddow, J. Hitschler, M. Junczys-Dowmunt, S. Läubli, A. V. Miceli Barone, J. Mokry, M. Nadejde, Nematus: A toolkit for neural machine translation, in: Proceedings of the Software Demonstrations of the 15th Conference of the European Chapter of the Association for Computational Linguistics, Association for Computational Linguistics, Valencia, Spain, 2017, pp. 65–68. URL: <https://www.aclweb.org/anthology/E17-3017>.
 - [28] M. Junczys-Dowmunt, T. Dwojak, H. Hoang, Is neural machine translation ready for deployment? A case study on 30 translation directions, in: International Workshop on Spoken Language

- Translation IWSLT, 2016.
- [29] Álvaro Peris, F. Casacuberta, NMT-Keras: A very flexible toolkit with a focus on interactive NMT and online learning, *The Prague Bulletin of Mathematical Linguistics* 111 (2018) 113–124.
 - [30] J. Helcl, J. Libovický, Neural Monkey: An open-source tool for sequence learning, *The Prague Bulletin of Mathematical Linguistics* (2017) 5–17.
 - [31] J. Zhang, Y. Ding, S. Shen, Y. Cheng, M. Sun, H. Luan, Y. Liu, THUMT: An open source toolkit for neural machine translation, *arXiv preprint arXiv:1706.06415* (2017).
 - [32] G. Neubig, M. Sperber, X. Wang, M. Felix, A. Matthews, S. Padmanabhan, Y. Qi, D. Sachan, P. Arthur, P. Godard, J. Hewitt, R. Riad, L. Wang, XNMT: The eXtensible neural machine translation toolkit, in: *Proceedings of the 13th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Papers)*, Association for Machine Translation in the Americas, Boston, MA, 2018, pp. 185–192. URL: <https://www.aclweb.org/anthology/W18-1818>.
 - [33] F. Stahlberg, E. Hasler, D. Saunders, B. Byrne, SGNMT – a flexible NMT decoding platform for quick prototyping of new models and search strategies, in: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 25–30. URL: <https://www.aclweb.org/anthology/D17-2005>. doi:10.18653/v1/D17-2005.
 - [34] F. Stahlberg, D. Saunders, G. Iglesias, B. Byrne, Why not be versatile? Applications of the SGNMT decoder for machine translation, in: *Proceedings of the 13th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Papers)*, Association for Machine Translation in the Americas, Boston, MA, 2018, pp. 208–216. URL: <https://www.aclweb.org/anthology/W18-1821>.
 - [35] X. Wang, M. Utiyama, E. Sumita, CytonMT: An efficient neural machine translation open-source toolkit implemented in C++, in: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 133–138. URL: <https://www.aclweb.org/anthology/D18-2023>.
 - [36] H. Xu, Q. Liu, Neutron: An implementation of the Transformer translation model and its variants, *arXiv preprint arXiv:1903.07402* (2019).
 - [37] J. R. Bellegarda, A latent semantic analysis framework for large-span language modeling, in: *Eurospeech*, 1997.
 - [38] R. Collobert, J. Weston, A unified architecture for natural language processing: Deep neural networks with multitask learning, in: *Proceedings of the 25th International Conference on Machine Learning, ICML '08*, ACM, New York, NY, USA, 2008, pp. 160–167. URL: <http://doi.acm.org/10.1145/1390156.1390177>. doi:10.1145/1390156.1390177.
 - [39] D. E. Rumelhart, G. E. Hinton, R. J. Williams, *Neurocomputing: Foundations of Research*, MIT Press, Cambridge, MA, USA, 1988, pp. 696–699. URL: <http://dl.acm.org/citation.cfm?id=65669.104451>.
 - [40] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, P. Kuksa, Natural language processing (almost) from scratch, *J. Mach. Learn. Res.* 12 (2011) 2493–2537.
 - [41] J. Pennington, R. Socher, C. D. Manning, GloVe: Global vectors for word representation, in: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Association for Computational Linguistics, Doha, Qatar, 2014, pp. 1532–1543. URL: <https://www.aclweb.org/anthology/D14-1162>. doi:10.3115/v1/D14-1162.
 - [42] T. Mikolov, K. Chen, G. Corrado, J. Dean, Efficient estimation of word representations in vector space, *CoRR* abs/1301.3781 (2013).
 - [43] T. Mikolov, Q. V. Le, I. Sutskever, Exploiting similarities among languages for machine translation, *arXiv preprint arXiv:1309.4168* (2013).
 - [44] S. Upadhyay, M. Faruqui, C. Dyer, D. Roth, Cross-lingual models of word embeddings: An empirical comparison, in: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Association for Computational Linguistics, Berlin, Germany, 2016, pp. 1661–1670. URL: <https://www.aclweb.org/anthology/P16-1157>. doi:10.18653/v1/P16-1157.
 - [45] M. Peters, W. Ammar, C. Bhagavatula, R. Power, Semi-supervised sequence tagging with bidirectional language models, in: *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Association for Computational Linguistics, 2017, pp. 1756–1765. URL: <http://aclweb.org/anthology/P17-1161>. doi:10.18653/v1/P17-1161.
 - [46] B. McCann, J. Bradbury, C. Xiong, R. Socher, Learned in translation: Contextualized word vectors, in: I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett (Eds.), *Advances in Neural Information Processing Sys-*

- tems 30, Curran Associates, Inc., 2017, pp. 6294–6305. URL: <http://papers.nips.cc/paper/7209-learned-in-translation-contextualized-word-vectors.pdf>.
- [47] M. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, L. Zettlemoyer, Deep contextualized word representations, in: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), Association for Computational Linguistics, 2018, pp. 2227–2237. URL: <http://aclweb.org/anthology/N18-1202>. doi:10.18653/v1/N18-1202.
 - [48] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, Improving language understanding with unsupervised learning, Technical Report, Technical report, OpenAI, 2018.
 - [49] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, BERT: Pre-training of deep bidirectional Transformers for language understanding, 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), Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 4171–4186. URL: <https://www.aclweb.org/anthology/N19-1423>.
 - [50] S. Bowman, E. Pavlick, E. Grave, B. van Durme, A. Wang, J. Hula, P. Xia, R. Pappagari, R. T. McCoy, R. Patel, et al., Looking for ELMo’s friends: Sentence-level pretraining beyond language modeling, arXiv preprint arXiv:1812.10860 (2018).
 - [51] Y. Goldberg, Assessing BERT’s syntactic abilities, arXiv preprint arXiv:1901.05287 (2019).
 - [52] H. Choi, K. Cho, Y. Bengio, Context-dependent word representation for neural machine translation, *Computer Speech & Language* 45 (2017) 149 – 160.
 - [53] J. B. Pollack, Recursive distributed representations, *Artificial Intelligence* 46 (1990) 77 – 105.
 - [54] R. Socher, J. Pennington, E. H. Huang, A. Y. Ng, C. D. Manning, Semi-supervised recursive Autoencoders for predicting sentiment distributions, in: Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Edinburgh, Scotland, UK., 2011, pp. 151–161. URL: <https://www.aclweb.org/anthology/D11-1014>.
 - [55] D. Wu, Stochastic inversion transduction grammars and bilingual parsing of parallel corpora, *Computational Linguistics* 23 (1997) 377–403.
 - [56] P. Li, Y. Liu, M. Sun, Recursive autoencoders for ITG-based translation, in: Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Seattle, Washington, USA, 2013, pp. 567–577. URL: <https://www.aclweb.org/anthology/D13-1054>.
 - [57] N. Kalchbrenner, E. Grefenstette, P. Blunsom, A convolutional neural network for modelling sentences, in: Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, 2014, pp. 655–665. URL: <http://aclweb.org/anthology/P14-1062>. doi:10.3115/v1/P14-1062.
 - [58] Y. Kim, Convolutional neural networks for sentence classification, in: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), Association for Computational Linguistics, 2014, pp. 1746–1751. URL: <http://aclweb.org/anthology/D14-1181>. doi:10.3115/v1/D14-1181.
 - [59] L. Mou, R. Men, G. Li, Y. Xu, L. Zhang, R. Yan, Z. Jin, Natural language inference by tree-based convolution and heuristic matching, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Association for Computational Linguistics, 2016, pp. 130–136. URL: <http://aclweb.org/anthology/P16-2022>. doi:10.18653/v1/P16-2022.
 - [60] C. dos Santos, M. Gatti, Deep convolutional neural networks for sentiment analysis of short texts, in: Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers, Dublin City University and Association for Computational Linguistics, 2014, pp. 69–78. URL: <http://aclweb.org/anthology/C14-1008>.
 - [61] M. J. Er, Y. Zhang, N. Wang, M. Pratama, Attention pooling-based convolutional neural network for sentence modelling, *Information Sciences* 373 (2016) 388 – 403.
 - [62] T. Shen, T. Zhou, G. Long, J. Jiang, S. Pan, C. Zhang, DiSAN: Directional self-attention network for RNN/CNN-free language understanding, in: Thirty-Second AAAI Conference on Artificial Intelligence, 2018.
 - [63] W. Wu, H. Wang, T. Liu, S. Ma, Phrase-level self-attention networks for universal sentence encoding, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 3729–3738. URL: <https://www.aclweb.org/anthology/D18-1408>.
 - [64] Q. Zhang, S. Liang, E. Yilmaz, Variational self-attention model for sentence representation, arXiv preprint arXiv:1812.11559 (2018).

- [65] L. Yu, C. d. M. d’Autume, C. Dyer, P. Blunsom, L. Kong, W. Ling, Sentence encoding with tree-constrained relation networks, arXiv preprint arXiv:1811.10475 (2018).
- [66] A. Santoro, D. Raposo, D. G. Barrett, M. Malinowski, R. Pascanu, P. Battaglia, T. Lillicrap, A simple neural network module for relational reasoning, in: I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett (Eds.), *Advances in Neural Information Processing Systems 30*, Curran Associates, Inc., 2017, pp. 4967–4976. URL: <http://papers.nips.cc/paper/7082-a-simple-neural-network-module-for-relational-reasoning.pdf>.
- [67] R. Palm, U. Paquet, O. Winther, Recurrent relational networks, in: S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, R. Garnett (Eds.), *Advances in Neural Information Processing Systems 31*, Curran Associates, Inc., 2018, pp. 3368–3378. URL: <http://papers.nips.cc/paper/7597-recurrent-relational-networks.pdf>.
- [68] A. Conneau, D. Kiela, H. Schwenk, L. Barrault, A. Bordes, Supervised learning of universal sentence representations from natural language inference data, in: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, 2017, pp. 670–680. URL: <http://aclweb.org/anthology/D17-1070>. doi:10.18653/v1/D17-1070.
- [69] A. Conneau, G. Kruszewski, G. Lample, L. Barrault, M. Baroni, What you can cram into a single vector: Probing sentence embeddings for linguistic properties, in: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Association for Computational Linguistics, 2018, pp. 2126–2136. URL: <http://aclweb.org/anthology/P18-1198>.
- [70] J. Wieting, D. Kiela, No training required: Exploring random encoders for sentence classification, arXiv preprint arXiv:1901.10444 (2019).
- [71] S. Hochreiter, J. Schmidhuber, Long short-term memory, *Neural Comput.* 9 (1997) 1735–1780.
- [72] S. Hochreiter, Y. Bengio, P. Frasconi, J. Schmidhuber, Gradient flow in recurrent nets: The difficulty of learning long-term dependencies, 2001.
- [73] M. Schuster, K. K. Paliwal, Bidirectional recurrent neural networks, *IEEE Transactions on Signal Processing* 45 (1997) 2673–2681.
- [74] P. Soutsov, S. Sarawagi, Length bias in encoder decoder models and a case for global conditioning, in: *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Austin, Texas, 2016, pp. 1516–1525. URL: <https://www.aclweb.org/anthology/D16-1158>. doi:10.18653/v1/D16-1158.
- [75] J. Pouget-Abadie, D. Bahdanau, B. van Merriënboer, K. Cho, Y. Bengio, Overcoming the curse of sentence length for neural machine translation using automatic segmentation, in: *Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation*, Association for Computational Linguistics, Doha, Qatar, 2014, pp. 78–85. URL: <https://www.aclweb.org/anthology/W14-4009>. doi:10.3115/v1/W14-4009.
- [76] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin, Attention is all you need, in: I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett (Eds.), *Advances in Neural Information Processing Systems 30*, Curran Associates, Inc., 2017, pp. 5998–6008. URL: <http://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf>.
- [77] M.-T. Luong, H. Pham, C. D. Manning, Effective approaches to attention-based neural machine translation, in: *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, 2015, pp. 1412–1421. URL: <http://aclweb.org/anthology/D15-1166>. doi:10.18653/v1/D15-1166.
- [78] H. Mino, M. Utiyama, E. Sumita, T. Tokunaga, Key-value attention mechanism for neural machine translation, in: *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, Asian Federation of Natural Language Processing, Taipei, Taiwan, 2017, pp. 290–295. URL: <https://www.aclweb.org/anthology/I17-2049>.
- [79] H. Larochelle, G. E. Hinton, Learning to combine foveal glimpses with a third-order Boltzmann machine, in: *Advances in neural information processing systems*, 2010, pp. 1243–1251.
- [80] J. L. Ba, V. Mnih, K. Kavukcuoglu, Multiple object recognition with visual attention, arXiv preprint arXiv:1412.7755 (2014).
- [81] V. Mnih, N. Heess, A. Graves, et al., Recurrent models of visual attention, in: *Advances in Neural Information Processing Systems*, 2014, pp. 2204–2212.
- [82] K. Xu, J. L. Ba, J. R. Kiros, K. Cho, A. Courville, R. Salakhutdinov, R. Zemel, Y. Bengio, Show, attend and tell: Neural image caption generation with visual attention, in: *International conference on machine learning*, 2015, pp. 2048–2057.
- [83] L. Yao, A. Torabi, K. Cho, N. Ballas, C. Pal, H. Larochelle, A. Courville, Describing videos by exploiting temporal structure, in: *Proceedings of the IEEE international conference on computer*

- p vision, 2015, pp. 4507–4515.
- [84] J. Chorowski, D. Bahdanau, K. Cho, Y. Bengio, End-to-end continuous speech recognition using attention-based recurrent NN: First results, in: NIPS 2014 Workshop on Deep Learning, December 2014, 2014.
 - [85] W. Chan, N. Jaitly, Q. V. Le, O. Vinyals, Listen, attend and spell: A neural network for large vocabulary conversational speech recognition, in: 2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2016, pp. 4960–4964. doi:10.1109/ICASSP.2016.7472621.
 - [86] L. Duong, A. Anastasopoulos, D. Chiang, S. Bird, T. Cohn, An attentional model for speech translation without transcription, in: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, San Diego, California, 2016, pp. 949–959. URL: <https://www.aclweb.org/anthology/N16-1109>. doi:10.18653/v1/N16-1109.
 - [87] S. K. Sønderby, C. K. Sønderby, H. Nielsen, O. Winther, Convolutional LSTM networks for subcellular localization of proteins, in: A.-H. Dediu, F. Hernández-Quiroz, C. Martín-Vide, D. A. Rosenblueth (Eds.), Algorithms for Computational Biology, Springer International Publishing, Cham, 2015, pp. 68–80.
 - [88] A. M. Rush, S. Chopra, J. Weston, A neural attention model for abstractive sentence summarization, in: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Lisbon, Portugal, 2015, pp. 379–389. URL: <https://www.aclweb.org/anthology/D15-1044>. doi:10.18653/v1/D15-1044.
 - [89] R. Sproat, N. Jaitly, RNN approaches to text normalization: A challenge, arXiv preprint arXiv:1611.00068 (2016).
 - [90] Z. Yuan, T. Briscoe, Grammatical error correction using neural machine translation, in: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, San Diego, California, 2016, pp. 380–386. URL: <https://www.aclweb.org/anthology/N16-1042>. doi:10.18653/v1/N16-1042.
 - [91] K. M. Hermann, T. Kocisky, E. Grefenstette, L. Espeholt, W. Kay, M. Suleyman, P. Blunsom, Teaching machines to read and comprehend, in: C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, R. Garnett (Eds.), Advances in Neural Information Processing Systems 28, Curran Associates, Inc., 2015, pp. 1693–1701. URL: <http://papers.nips.cc/paper/5945-teaching-machines-to-read-and-comprehend.pdf>.
 - [92] Z. Yang, X. He, J. Gao, L. Deng, A. Smola, Stacked attention networks for image question answering, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 21–29.
 - [93] S. Sukhbaatar, A. Szlam, J. Weston, R. Fergus, End-to-end memory networks, in: C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, R. Garnett (Eds.), Advances in Neural Information Processing Systems 28, Curran Associates, Inc., 2015, pp. 2440–2448. URL: <http://papers.nips.cc/paper/5846-end-to-end-memory-networks.pdf>.
 - [94] L. Dong, M. Lapata, Language to logical form with neural attention, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Berlin, Germany, 2016, pp. 33–43. URL: <https://www.aclweb.org/anthology/P16-1004>. doi:10.18653/v1/P16-1004.
 - [95] J. Im, S. Cho, Distance-based self-attention network for natural language inference, arXiv preprint arXiv:1712.02047 (2017).
 - [96] Y. Liu, C. Sun, L. Lin, X. Wang, Learning natural language inference using bidirectional LSTM model and inner-attention, arXiv preprint arXiv:1605.09090 (2016).
 - [97] H. Adel, H. Schütze, Exploring different dimensions of attention for uncertainty detection, in: Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers, Association for Computational Linguistics, Valencia, Spain, 2017, pp. 22–34. URL: <https://www.aclweb.org/anthology/E17-1003>.
 - [98] C.-Y. Lee, S. Osindero, Recursive recurrent nets with attention modeling for OCR in the wild, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016.
 - [99] L. Shang, Z. Lu, H. Li, Neural responding machine for short-text conversation, in: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Association for Computational Linguistics, Beijing, China, 2015, pp. 1577–1586. URL: <https://www.aclweb.org/anthology/P15-1152>. doi:10.3115/v1/P15-1152.
 - [100] Basho, J. Reichhold, Basho: the complete haiku, Kodansha International, 2013.

- [101] F. Bastien, P. Lamblin, R. Pascanu, J. Bergstra, I. Goodfellow, A. Bergeron, N. Bouchard, D. Warde-Farley, Y. Bengio, Theano: New features and speed improvements, in: NIPS, 2012.
- [102] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. G. Murray, B. Steiner, P. Tucker, V. Vasudevan, P. Warden, M. Wicke, Y. Yu, X. Zheng, Tensorflow: A system for large-scale machine learning, in: 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16), USENIX Association, Savannah, GA, 2016, pp. 265–283. URL: <https://www.usenix.org/conference/osdi16/technical-sessions/presentation/abadi>.
- [103] I. Goodfellow, D. Warde-farley, M. Mirza, A. Courville, Y. Bengio, Maxout networks, in: ICML, 2013, pp. 1319–1327.
- [104] A. Waibel, T. Hanazawa, G. E. Hinton, K. Shikano, K. J. Lang, Phoneme recognition using time-delay neural networks, IEEE Transactions on Acoustics, Speech, and Signal Processing 37 (1989) 328–339.
- [105] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. E. Hubbard, L. D. Jackel, Backpropagation applied to handwritten zip code recognition, Neural Comput. 1 (1989) 541–551.
- [106] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. E. Hubbard, L. D. Jackel, Handwritten digit recognition with a back-propagation network, in: D. S. Touretzky (Ed.), Advances in Neural Information Processing Systems 2, Morgan-Kaufmann, 1990, pp. 396–404. URL: <http://papers.nips.cc/paper/293-handwritten-digit-recognition-with-a-back-propagation-network.pdf>.
- [107] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, Gradient-based learning applied to document recognition, Proceedings of the IEEE 86 (1998) 2278–2324.
- [108] F. Chollet, Xception: Deep learning with depthwise separable convolutions, in: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, 2017, pp. 1800–1807.
- [109] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, H. Adam, Mobilenets: Efficient convolutional neural networks for mobile vision applications, arXiv preprint arXiv:1704.04861 (2017).
- [110] J. Gehring, M. Auli, D. Grangier, D. Yarats, Y. N. Dauphin, Convolutional sequence to sequence learning, in: Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML’17, JMLR.org, 2017, pp. 1243–1252. URL: <http://dl.acm.org/citation.cfm?id=3305381.3305510>.
- [111] N. Kalchbrenner, L. Espeholt, K. Simonyan, A. van den Oord, A. Graves, K. Kavukcuoglu, Neural machine translation in linear time, arXiv preprint arXiv:1610.10099 (2016).
- [112] J. Gehring, M. Auli, D. Grangier, Y. N. Dauphin, A convolutional encoder model for neural machine translation, in: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 123–135. URL: <https://www.aclweb.org/anthology/P17-1012>. doi:10.18653/v1/P17-1012.
- [113] L. Kaiser, A. N. Gomez, F. Chollet, Depthwise separable convolutions for neural machine translation, arXiv preprint arXiv:1706.03059 (2017).
- [114] F. Wu, A. Fan, A. Baevski, Y. N. Dauphin, M. Auli, Pay less attention with lightweight and dynamic convolutions, in: ICLR, 2019.
- [115] A. van den Oord, N. Kalchbrenner, K. Kavukcuoglu, Pixel recurrent neural networks, in: Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48, ICML’16, JMLR.org, 2016, pp. 1747–1756. URL: <http://dl.acm.org/citation.cfm?id=3045390.3045575>.
- [116] J. Cheng, L. Dong, M. Lapata, Long short-term memory-networks for machine reading, in: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Austin, Texas, 2016, pp. 551–561. URL: <https://www.aclweb.org/anthology/D16-1053>. doi:10.18653/v1/D16-1053.
- [117] A. Parikh, O. Täckström, D. Das, J. Uszkoreit, A decomposable attention model for natural language inference, in: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Austin, Texas, 2016, pp. 2249–2255. URL: <https://www.aclweb.org/anthology/D16-1244>. doi:10.18653/v1/D16-1244.
- [118] T. Shen, T. Zhou, G. Long, J. Jiang, S. Wang, C. Zhang, Reinforced self-attention network: A hybrid of hard and soft attention for sequence modeling, in: Proceedings of the 27th International Joint Conference on Artificial Intelligence, IJCAI’18, AAAI Press, 2018, pp. 4345–4352. URL: <http://dl.acm.org/citation.cfm?id=3304222.3304374>.
- [119] R. Paulus, C. Xiong, R. Socher, A deep reinforced model for abstractive summarization, arXiv

- preprint arXiv:1705.04304 (2017).
- [120] G. Daniil, P. Kalaidin, V. Malykh, Self-attentive model for headline generation, arXiv preprint arXiv:1901.07786 (2019).
 - [121] Z. Lin, M. Feng, C. N. d. Santos, M. Yu, B. Xiang, B. Zhou, Y. Bengio, A structured self-attentive sentence embedding, arXiv preprint arXiv:1703.03130 (2017).
 - [122] M. Hu, Y. Peng, Z. Huang, X. Qiu, F. Wei, M. Zhou, Reinforced mnemonic reader for machine reading comprehension, in: Proceedings of the 27th International Joint Conference on Artificial Intelligence, IJCAI'18, AAAI Press, 2018, pp. 4099–4106. URL: <http://dl.acm.org/citation.cfm?id=3304222.3304340>.
 - [123] G. Tang, M. Müller, A. Rios, R. Sennrich, Why self-attention? A targeted evaluation of neural machine translation architectures, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 4263–4272. URL: <https://www.aclweb.org/anthology/D18-1458>.
 - [124] M. X. Chen, O. Firat, A. Bapna, M. Johnson, W. Macherey, G. Foster, L. Jones, M. Schuster, N. Shazeer, N. Parmar, A. Vaswani, J. Uszkoreit, L. Kaiser, Z. Chen, Y. Wu, M. Hughes, The best of both worlds: Combining recent advances in neural machine translation, in: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 76–86. URL: <https://www.aclweb.org/anthology/P18-1008>.
 - [125] P. Shaw, J. Uszkoreit, A. Vaswani, Self-attention with relative position representations, in: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), Association for Computational Linguistics, New Orleans, Louisiana, 2018, pp. 464–468. URL: <https://www.aclweb.org/anthology/N18-2074>. doi:10.18653/v1/N18-2074.
 - [126] F. Stahlberg, A. de Gispert, B. Byrne, The University of Cambridge’s machine translation systems for WMT18, in: Proceedings of the Third Conference on Machine Translation: Shared Task Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 504–512. URL: <https://www.aclweb.org/anthology/W18-6427>.
 - [127] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov, Dropout: A simple way to prevent neural networks from overfitting, The Journal of Machine Learning Research 15 (2014) 1929–1958.
 - [128] S. Ioffe, C. Szegedy, Batch normalization: Accelerating deep network training by reducing internal covariate shift, in: Proceedings of the 32Nd International Conference on International Conference on Machine Learning - Volume 37, ICML’15, JMLR.org, 2015, pp. 448–456. URL: <http://dl.acm.org/citation.cfm?id=3045118.3045167>.
 - [129] J. L. Ba, J. R. Kiros, G. E. Hinton, Layer normalization, arXiv preprint arXiv:1607.06450 (2016).
 - [130] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.
 - [131] L. Miculicich, N. Pappas, D. Ram, A. Popescu-Belis, Self-attentive residual decoder for neural machine translation, in: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), Association for Computational Linguistics, New Orleans, Louisiana, 2018, pp. 1366–1379. URL: <https://www.aclweb.org/anthology/N18-1124>. doi:10.18653/v1/N18-1124.
 - [132] C. Wang, M. Li, A. Smola, Language models with Transformers, arXiv preprint arXiv:1904.09408 (2019).
 - [133] L. M. Werlen, N. Pappas, D. Ram, A. Popescu-Belis, Global-context neural machine translation through target-side attentive residual connections, researchgate.net (2018).
 - [134] J. Hao, X. Wang, B. Yang, L. Wang, J. Zhang, Z. Tu, Modeling recurrence for Transformer, arXiv preprint arXiv:1904.03092 (2019).
 - [135] J. Lin, X. Sun, X. Ren, S. Ma, J. Su, Q. Su, Deconvolution-based global decoding for neural machine translation, in: Proceedings of the 27th International Conference on Computational Linguistics, Association for Computational Linguistics, Santa Fe, New Mexico, USA, 2018, pp. 3260–3271. URL: <https://www.aclweb.org/anthology/C18-1276>.
 - [136] F. Stahlberg, B. Byrne, On NMT search errors and model errors: Cat got your tongue?, in: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Hong Kong, 2019.
 - [137] T. G. Dietterich, Ensemble methods in machine learning, in: International workshop on multiple classifier systems, Springer Berlin Heidelberg, Berlin, Heidelberg, 2000, pp. 1–15.
 - [138] L. K. Hansen, P. Salamon, Neural network ensembles, IEEE transactions on pattern analysis and

- machine intelligence 12 (1990) 993–1001.
- [139] G. Neubig, Lexicons and minimum risk training for neural machine translation: NAIST-CMU at WAT2016, in: Proceedings of the 3rd Workshop on Asian Translation (WAT2016), The COLING 2016 Organizing Committee, Osaka, Japan, 2016, pp. 119–125. URL: <https://www.aclweb.org/anthology/W16-4610>.
 - [140] R. Sennrich, B. Haddow, A. Birch, Edinburgh neural machine translation systems for WMT 16, in: Proceedings of the First Conference on Machine Translation, Association for Computational Linguistics, Berlin, Germany, 2016, pp. 371–376. URL: <https://www.aclweb.org/anthology/W16-2323>. doi:10.18653/v1/W16-2323.
 - [141] F. Cromieres, C. Chu, T. Nakazawa, S. Kurohashi, Kyoto university participation to WAT 2016, in: Proceedings of the 3rd Workshop on Asian Translation (WAT2016), The COLING 2016 Organizing Committee, Osaka, Japan, 2016, pp. 166–174. URL: <https://www.aclweb.org/anthology/W16-4616>.
 - [142] N. Durrani, F. Dalvi, H. Sajjad, S. Vogel, QCRI machine translation systems for IWSLT 16, in: International Workshop on Spoken Language Translation. Seattle, WA, USA, 2016.
 - [143] O. Bojar, R. Chatterjee, C. Federmann, Y. Graham, B. Haddow, M. Huck, A. Jimeno Yepes, P. Koehn, V. Logacheva, C. Monz, M. Negri, A. Neveol, M. Neves, M. Popel, M. Post, R. Rubino, C. Scarton, L. Specia, M. Turchi, K. Verspoor, M. Zampieri, Findings of the 2016 conference on machine translation, in: Proceedings of the First Conference on Machine Translation, Association for Computational Linguistics, Berlin, Germany, 2016, pp. 131–198. URL: <https://www.aclweb.org/anthology/W16-2301>. doi:10.18653/v1/W16-2301.
 - [144] O. Bojar, R. Chatterjee, C. Federmann, Y. Graham, B. Haddow, S. Huang, M. Huck, P. Koehn, Q. Liu, V. Logacheva, C. Monz, M. Negri, M. Post, R. Rubino, L. Specia, M. Turchi, Findings of the 2017 conference on machine translation (WMT17), in: Proceedings of the Second Conference on Machine Translation, Association for Computational Linguistics, 2017, pp. 169–214. URL: <http://aclweb.org/anthology/W17-4717>. doi:10.18653/v1/W17-4717.
 - [145] O. Bojar, C. Federmann, M. Fishel, Y. Graham, B. Haddow, P. Koehn, C. Monz, Findings of the 2018 conference on machine translation (WMT18), in: Proceedings of the Third Conference on Machine Translation: Shared Task Papers, Association for Computational Linguistics, 2018, pp. 272–303. URL: <http://aclweb.org/anthology/W18-6401>.
 - [146] O. Bojar, et al., Findings of the 2019 conference on machine translation (WMT19), in: Proceedings of the Fourth Conference on Machine Translation: Shared Task Papers, Association for Computational Linguistics, 2019.
 - [147] R. Sennrich, A. Birch, A. Currey, U. Germann, B. Haddow, K. Heafield, A. V. Miceli Barone, P. Williams, The University of Edinburgh’s neural MT systems for WMT17, in: Proceedings of the Second Conference on Machine Translation, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 389–399. URL: <https://www.aclweb.org/anthology/W17-4739>. doi:10.18653/v1/W17-4739.
 - [148] Y. Wang, S. Cheng, L. Jiang, J. Yang, W. Chen, M. Li, L. Shi, Y. Wang, H. Yang, Sogou neural machine translation systems for WMT17, in: Proceedings of the Second Conference on Machine Translation, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 410–415. URL: <https://www.aclweb.org/anthology/W17-4742>. doi:10.18653/v1/W17-4742.
 - [149] M. Junczys-Dowmunt, Microsoft’s submission to the WMT2018 news translation task: How I learned to stop worrying and love the data, in: Proceedings of the Third Conference on Machine Translation: Shared Task Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 425–430. URL: <https://www.aclweb.org/anthology/W18-6415>.
 - [150] M. Wang, L. Gong, W. Zhu, J. Xie, C. Bian, Tencent neural machine translation systems for WMT18, in: Proceedings of the Third Conference on Machine Translation: Shared Task Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 522–527. URL: <https://www.aclweb.org/anthology/W18-6429>.
 - [151] F. Stahlberg, B. Byrne, Unfolding and shrinking neural machine translation ensembles, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 1946–1956. URL: <https://www.aclweb.org/anthology/D17-1208>. doi:10.18653/v1/D17-1208.
 - [152] M. Freitag, Y. Al-Onaizan, Fast domain adaptation for neural machine translation, arXiv preprint arXiv:1612.06897 (2016).
 - [153] X. He, G. Haffari, M. Norouzi, Sequence to sequence mixture model for diverse machine translation, in: Proceedings of the 22nd Conference on Computational Natural Language Learning, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 583–592. URL: <https://www.aclweb.org/anthology/W18-4401>.

- org/anthology/K18-1056.
- [154] O. Firat, B. Sankaran, Y. Al-Onaizan, F. T. Yarman Vural, K. Cho, Zero-resource translation with multi-lingual neural machine translation, in: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Austin, Texas, 2016, pp. 268–277. URL: <https://www.aclweb.org/anthology/D16-1026>. doi:10.18653/v1/D16-1026.
 - [155] B. Zoph, K. Knight, Multi-source neural translation, in: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, San Diego, California, 2016, pp. 30–34. URL: <https://www.aclweb.org/anthology/N16-1004>. doi:10.18653/v1/N16-1004.
 - [156] L. Rokach, Ensemble-based classifiers, *Artificial Intelligence Review* 33 (2010) 1–39.
 - [157] R. Sennrich, B. Haddow, A. Birch, Neural machine translation of rare words with subword units, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Berlin, Germany, 2016, pp. 1715–1725. URL: <https://www.aclweb.org/anthology/P16-1162>. doi:10.18653/v1/P16-1162.
 - [158] I. Goodfellow, Y. Bengio, A. Courville, *Deep Learning*, MIT Press, 2016. <http://www.deeplearningbook.org>.
 - [159] M. Junczys-Dowmunt, T. Dwojak, R. Sennrich, The AMU-UEDIN submission to the WMT16 news translation task: Attention-based NMT models as feature functions in phrase-based SMT, in: Proceedings of the First Conference on Machine Translation, Association for Computational Linguistics, Berlin, Germany, 2016, pp. 319–325. URL: <https://www.aclweb.org/anthology/W16-2316>. doi:10.18653/v1/W16-2316.
 - [160] Y. Liu, L. Zhou, Y. Wang, Y. Zhao, J. Zhang, C. Zong, A comparable study on model averaging, ensembling and reranking in NMT, in: M. Zhang, V. Ng, D. Zhao, S. Li, H. Zan (Eds.), *Natural Language Processing and Chinese Computing*, Springer International Publishing, Cham, 2018, pp. 299–308.
 - [161] M. Popel, O. Bojar, Training tips for the Transformer model, *The Prague Bulletin of Mathematical Linguistics* 110 (2018) 43–70.
 - [162] L. Liu, M. Utiyama, A. Finch, E. Sumita, Agreement on target-bidirectional neural machine translation, in: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, San Diego, California, 2016, pp. 411–416. URL: <https://www.aclweb.org/anthology/N16-1046>. doi:10.18653/v1/N16-1046.
 - [163] Z. Zhang, S. Wu, S. Liu, M. Li, M. Zhou, E. Chen, Regularizing neural machine translation by target-bidirectional agreement, *arXiv preprint arXiv:1808.04064* (2018).
 - [164] Z. Yang, L. Chen, M. Le Nguyen, Regularizing forward and backward decoding to improve neural machine translation, in: 2018 10th International Conference on Knowledge and Systems Engineering (KSE), 2018, pp. 73–78. doi:10.1109/KSE.2018.8573433.
 - [165] S. Mehri, L. Sigal, Middle-out decoding, in: S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, R. Garnett (Eds.), *Advances in Neural Information Processing Systems* 31, Curran Associates, Inc., 2018, pp. 5518–5529. URL: <http://papers.nips.cc/paper/7796-middle-out-decoding.pdf>.
 - [166] S. Welleck, K. Brantley, H. Daumé III, K. Cho, Non-monotonic sequential text generation, *arXiv preprint arXiv:1902.02192* (2019).
 - [167] J. Gu, Q. Liu, K. Cho, Insertion-based decoding with automatically inferred generation order, *arXiv preprint arXiv:1902.01370* (2019).
 - [168] M. Stern, W. Chan, J. R. Kiros, J. Uszkoreit, Insertion Transformer: Flexible sequence generation via insertion operations, *arXiv preprint arXiv:1902.03249* (2019).
 - [169] R. Östling, J. Tiedemann, Neural machine translation for low-resource languages, *arXiv preprint arXiv:1708.05729* (2017).
 - [170] J. Gu, C. Wang, J. Zhao, Levenshtein Transformer, *arXiv preprint arXiv:1905.11006* (2019).
 - [171] F. Stahlberg, D. Saunders, B. Byrne, An operation sequence model for explainable neural machine translation, in: Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 175–186. URL: <https://www.aclweb.org/anthology/W18-5420>.
 - [172] A. Li, S. Zhang, D. Wang, T. F. Zheng, Enhanced neural machine translation by learning from draft, in: 2017 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC), 2017, pp. 1583–1587. doi:10.1109/APSIPA.2017.8282276.
 - [173] X. Zhang, J. Su, Y. Qin, Y. Liu, R. Ji, H. Wang, Asynchronous bidirectional decoding for neural

- machine translation, in: Thirty-Second AAAI Conference on Artificial Intelligence, 2018.
- [174] A. ElMaghraby, A. Rafea, Enhancing translation from English to Arabic using two-phase decoder translation, in: K. Arai, S. Kapoor, R. Bhatia (Eds.), *Intelligent Systems and Applications*, Springer International Publishing, Cham, 2019, pp. 539–549.
 - [175] X. Geng, X. Feng, B. Qin, T. Liu, Adaptive multi-pass decoder for neural machine translation, in: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 523–532. URL: <https://www.aclweb.org/anthology/D18-1048>.
 - [176] D. He, H. Lu, Y. Xia, T. Qin, L. Wang, T.-Y. Liu, Decoding with value networks for neural machine translation, in: I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett (Eds.), *Advances in Neural Information Processing Systems* 30, Curran Associates, Inc., 2017, pp. 178–187. URL: <http://papers.nips.cc/paper/6622-decoding-with-value-networks-for-neural-machine-translation.pdf>.
 - [177] V. C. D. Hoang, G. Haffari, T. Cohn, Towards decoding as continuous optimisation in neural machine translation, in: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 146–156. URL: <https://www.aclweb.org/anthology/D17-1014>. doi:10.18653/v1/D17-1014.
 - [178] M. Stern, N. Shazeer, J. Uszkoreit, Blockwise parallel decoding for deep autoregressive models, in: S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, R. Garnett (Eds.), *Advances in Neural Information Processing Systems* 31, Curran Associates, Inc., 2018, pp. 10086–10095. URL: <http://papers.nips.cc/paper/8212-blockwise-parallel-decoding-for-deep-autoregressive-models.pdf>.
 - [179] L. Kaiser, A. Roy, A. Vaswani, N. Pamar, S. Bengio, J. Uszkoreit, N. Shazeer, Fast decoding in sequence models using discrete latent variables, arXiv preprint arXiv:1803.03382 (2018).
 - [180] M. A. Di Gangi, M. Federico, Deep neural machine translation with weakly-recurrent units, in: *Proceedings of the 21st Annual Conference of the European Association for Machine Translation*: 28-30 May 2018, Universitat d’Alacant, Alacant, Spain, European Association for Machine Translation, 2018, pp. 119–128.
 - [181] P. Koehn, *Statistical Machine Translation*, 1st ed., Cambridge University Press, New York, NY, USA, 2010.
 - [182] Z. Zhang, R. Wang, M. Utiyama, E. Sumita, H. Zhao, Exploring recombination for efficient decoding of neural machine translation, in: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 4785–4790. URL: <https://www.aclweb.org/anthology/D18-1511>.
 - [183] X. Liu, Y. Wang, X. Chen, M. J. Gales, P. C. Woodland, Efficient lattice rescoring using recurrent neural network language models, in: 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2014, pp. 4908–4912. doi:10.1109/ICASSP.2014.6854535.
 - [184] G. Lecorvé, P. Motlicek, Conversion of recurrent neural network language models to weighted finite state transducers for automatic speech recognition, in: *Thirteenth Annual Conference of the International Speech Communication Association*, 2012.
 - [185] M. Freitag, Y. Al-Onaizan, Beam search strategies for neural machine translation, in: *Proceedings of the First Workshop on Neural Machine Translation*, Association for Computational Linguistics, Vancouver, 2017, pp. 56–60. URL: <https://www.aclweb.org/anthology/W17-3207>. doi:10.18653/v1/W17-3207.
 - [186] K. Goyal, G. Neubig, C. Dyer, T. Berg-Kirkpatrick, A continuous relaxation of beam search for end-to-end training of neural sequence models, in: *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
 - [187] S. Wiseman, A. M. Rush, Sequence-to-sequence learning as beam-search optimization, in: *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Austin, Texas, 2016, pp. 1296–1306. URL: <https://www.aclweb.org/anthology/D16-1137>. doi:10.18653/v1/D16-1137.
 - [188] R. Collobert, A. Hannun, G. Synnaeve, A fully differentiable beam search decoder, arXiv preprint arXiv:1902.06022 (2019).
 - [189] J. Gu, K. Cho, V. O. Li, Trainable greedy decoding for neural machine translation, in: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 1968–1978. URL: <https://www.aclweb.org/anthology/D17-1210>. doi:10.18653/v1/D17-1210.
 - [190] Y. Kim, A. M. Rush, Sequence-level knowledge distillation, in: *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics,

- Austin, Texas, 2016, pp. 1317–1327. URL: <https://www.aclweb.org/anthology/D16-1139>. doi:10.18653/v1/D16-1139.
- [191] Y. Chen, V. O. Li, K. Cho, S. Bowman, A stable and effective learning strategy for trainable greedy decoding, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 380–390. URL: <https://www.aclweb.org/anthology/D18-1035>.
 - [192] J. Li, D. Jurafsky, Mutual information and diverse decoding improve neural machine translation, arXiv preprint arXiv:1601.00372 (2016).
 - [193] J. Li, M. Galley, C. Brockett, J. Gao, B. Dolan, A diversity-promoting objective function for neural conversation models, in: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, San Diego, California, 2016, pp. 110–119. URL: <https://www.aclweb.org/anthology/N16-1014>. doi:10.18653/v1/N16-1014.
 - [194] K. Gimpel, D. Batra, C. Dyer, G. Shakhnarovich, A systematic exploration of diversity in machine translation, in: Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Seattle, Washington, USA, 2013, pp. 1100–1111. URL: <https://www.aclweb.org/anthology/D13-1111>.
 - [195] V. Goel, S. Kumar, B. Byrne, Segmental minimum Bayes-risk ASR voting strategies, in: Interspeech, 2000, pp. 139–142.
 - [196] S. Kumar, B. Byrne, Minimum Bayes-risk decoding for statistical machine translation, in: HLT-NAACL 2004: Main Proceedings, Association for Computational Linguistics, Boston, Massachusetts, USA, 2004, pp. 169–176. URL: <https://www.aclweb.org/anthology/N04-1022>.
 - [197] R. Tromble, S. Kumar, F. J. Och, W. Macherey, Lattice Minimum Bayes-Risk decoding for statistical machine translation, in: Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Honolulu, Hawaii, 2008, pp. 620–629. URL: <https://www.aclweb.org/anthology/D08-1065>.
 - [198] G. Iglesias, W. Tambellini, A. de Gispert, E. Hasler, B. Byrne, Accelerating NMT batched beam decoding with LMBR posteriors for deployment, in: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 3 (Industry Papers), Association for Computational Linguistics, New Orleans - Louisiana, 2018, pp. 106–113. URL: <https://www.aclweb.org/anthology/N18-3013>. doi:10.18653/v1/N18-3013.
 - [199] F. Stahlberg, A. de Gispert, E. Hasler, B. Byrne, Neural machine translation by minimising the Bayes-risk with respect to syntactic translation lattices, in: Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers, Association for Computational Linguistics, Valencia, Spain, 2017, pp. 362–368. URL: <https://www.aclweb.org/anthology/E17-2058>.
 - [200] K. Cho, Noisy parallel approximate decoding for conditional recurrent language model, arXiv preprint arXiv:1605.03835 (2016).
 - [201] Y. Bengio, G. Mesnil, Y. N. Dauphin, S. Rifai, Better mixing via deep representations, in: S. Dasgupta, D. McAllester (Eds.), Proceedings of the 30th International Conference on Machine Learning, volume 28 of *Proceedings of Machine Learning Research*, PMLR, Atlanta, Georgia, USA, 2013, pp. 552–560. URL: <http://proceedings.mlr.press/v28/bengio13.html>.
 - [202] A. K. Vijayakumar, M. Cogswell, R. R. Selvaraju, Q. Sun, S. Lee, D. Crandall, D. Batra, Diverse beam search: Decoding diverse solutions from neural sequence models, arXiv preprint arXiv:1610.02424 (2016).
 - [203] Y. Park, H. Na, H. Lee, J. Lee, I. Song, An effective diverse decoding scheme for robust synonymous sentence translation, AMTA 2016, Vol. (2016) 53.
 - [204] M. Müller, T. S. Nguyen, J. Niehues, E. Cho, B. Krüger, T.-L. Ha, K. Kilgour, M. Sperber, M. Mediani, S. Stüker, A. Waibel, Lecture translator - speech translation framework for simultaneous lecture translation, in: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations, Association for Computational Linguistics, San Diego, California, 2016, pp. 82–86. URL: <https://www.aclweb.org/anthology/N16-3017>. doi:10.18653/v1/N16-3017.
 - [205] C. Fügen, A. Waibel, M. Kolss, Simultaneous translation of lectures and speeches, Machine translation 21 (2007) 209–252.
 - [206] W. D. Lewis, Skype translator: Breaking down language and hearing barriers, Translating and the Computer (TC37) 10 (2015) 125–149.
 - [207] T. Mieno, G. Neubig, S. Sakti, T. Toda, S. Nakamura, Speed or accuracy? A study in evaluation

- of simultaneous speech translation, in: Sixteenth Annual Conference of the International Speech Communication Association, 2015.
- [208] K. Cho, M. Esipova, Can neural machine translation do simultaneous translation?, arXiv preprint arXiv:1606.02012 (2016).
 - [209] A. Grissom II, H. He, J. Boyd-Graber, J. Morgan, H. Daumé III, Don’t until the final verb wait: Reinforcement learning for simultaneous machine translation, in: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), Association for Computational Linguistics, Doha, Qatar, 2014, pp. 1342–1352. URL: <https://www.aclweb.org/anthology/D14-1140>. doi:10.3115/v1/D14-1140.
 - [210] J. Gu, G. Neubig, K. Cho, V. O. Li, Learning to translate in real-time with neural machine translation, in: Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers, Association for Computational Linguistics, Valencia, Spain, 2017, pp. 1053–1062. URL: <https://www.aclweb.org/anthology/E17-1099>.
 - [211] M. Ma, L. Huang, H. Xiong, K. Liu, C. Zhang, Z. He, H. Liu, X. Li, H. Wang, Stacl: Simultaneous translation with integrated anticipation and controllable latency, arXiv preprint arXiv:1810.08398 (2018).
 - [212] M. Paulik, A. Waibel, Automatic translation from parallel speech: Simultaneous interpretation as mt training data, in: 2009 IEEE Workshop on Automatic Speech Recognition Understanding, 2009, pp. 496–501. doi:10.1109/ASRU.2009.5372880.
 - [213] M. Paulik, A. Waibel, Training speech translation from audio recordings of interpreter-mediated communication, Computer Speech & Language 27 (2013) 455 – 474. Special Issue on Speech-speech translation.
 - [214] H. He, J. Boyd-Graber, H. Daumé III, Interpretese vs. Translationese: The uniqueness of human strategies in simultaneous interpretation, in: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, San Diego, California, 2016, pp. 971–976. URL: <https://www.aclweb.org/anthology/N16-1111>. doi:10.18653/v1/N16-1111.
 - [215] K. Heafield, I. Pouzyrevsky, J. H. Clark, P. Koehn, Scalable modified Kneser-Ney language model estimation, in: Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Association for Computational Linguistics, Sofia, Bulgaria, 2013, pp. 690–696. URL: <https://www.aclweb.org/anthology/P13-2121>.
 - [216] J. Lin, C. Dyer, Data-intensive text processing with MapReduce, in: NAACL HLT 2010 Tutorial Abstracts, Association for Computational Linguistics, Los Angeles, California, 2010, pp. 1–2. URL: <https://www.aclweb.org/anthology/N10-4001>.
 - [217] D. Chiang, Hierarchical phrase-based translation, American Journal of Computational Linguistics 33 (2007) 201–228.
 - [218] G. K. Zipf, The psychology of language, in: Encyclopedia of psychology, Philosophical Library, 1946, pp. 332–341.
 - [219] S. Jean, K. Cho, R. Memisevic, Y. Bengio, On using very large target vocabulary for neural machine translation, in: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Association for Computational Linguistics, Beijing, China, 2015, pp. 1–10. URL: <https://www.aclweb.org/anthology/P15-1001>. doi:10.3115/v1/P15-1001.
 - [220] M.-T. Luong, I. Sutskever, Q. V. Le, O. Vinyals, W. Zaremba, Addressing the rare word problem in neural machine translation, in: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Association for Computational Linguistics, Beijing, China, 2015, pp. 11–19. URL: <https://www.aclweb.org/anthology/P15-1002>. doi:10.3115/v1/P15-1002.
 - [221] Q. V. Le, M.-T. Luong, I. Sutskever, O. Vinyals, W. Zaremba, Neural machine translation systems with rare word processing, 2016. US Patent App. 14/921,925.
 - [222] C. Gulcehre, S. Ahn, R. Nallapati, B. Zhou, Y. Bengio, Pointing the unknown words, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Berlin, Germany, 2016, pp. 140–149. URL: <https://www.aclweb.org/anthology/P16-1014>. doi:10.18653/v1/P16-1014.
 - [223] X. Li, J. Zhang, C. Zong, Towards zero unknown word in neural machine translation, in: Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI’16, AAAI Press, 2016, pp. 2852–2858. URL: <http://dl.acm.org/citation.cfm?id=3060832.3061020>.
 - [224] F. Li, D. Quan, W. Qiang, X. Tong, J. Zhu, Handling many-to-one unk translation for neural machine translation, in: Machine Translation: 13th China Workshop, CWMT 2017, Revised

- Selected Papers, Springer, 2017, pp. 102–111.
- [225] S. Li, J. Xu, Y. Zhang, Y. Chen, A method of unknown words processing for neural machine translation using HowNet, in: Machine Translation: 13th China Workshop, CWMT 2017, Revised Selected Papers, Springer, 2017, pp. 20–29.
 - [226] G. Miao, J. Xu, Y. Li, S. Li, Y. Chen, An unknown word processing method in NMT by integrating syntactic structure and semantic concept, in: Machine Translation: 13th China Workshop, CWMT 2017, Revised Selected Papers, Springer, 2017, pp. 43–54.
 - [227] T. Q. Nguyen, D. Chiang, Improving lexical choice in neural machine translation, in: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), Association for Computational Linguistics, New Orleans, Louisiana, 2018, pp. 334–343. URL: <https://www.aclweb.org/anthology/N18-1031>. doi:10.18653/v1/N18-1031.
 - [228] J. Andreas, D. Klein, When and why are log-linear models self-normalizing?, in: Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, Denver, Colorado, 2015, pp. 244–249. URL: <https://www.aclweb.org/anthology/N15-1027>. doi:10.3115/v1/N15-1027.
 - [229] M. Gutmann, A. Hyvärinen, Noise-contrastive estimation: A new estimation principle for unnormalized statistical models, in: Y. W. Teh, M. Titterton (Eds.), Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, volume 9 of *Proceedings of Machine Learning Research*, PMLR, Chia Laguna Resort, Sardinia, Italy, 2010, pp. 297–304. URL: <http://proceedings.mlr.press/v9/gutmann10a.html>.
 - [230] C. Dyer, Notes on noise contrastive estimation and negative sampling, arXiv preprint arXiv:1410.8251 (2014).
 - [231] A. Mnih, K. Kavukcuoglu, Learning word embeddings efficiently with noise-contrastive estimation, in: C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, K. Q. Weinberger (Eds.), Advances in Neural Information Processing Systems 26, Curran Associates, Inc., 2013, pp. 2265–2273. URL: <http://papers.nips.cc/paper/5165-learning-word-embeddings-efficiently-with-noise-contrastive-estimation.pdf>.
 - [232] A. Mnih, Y. W. Teh, A fast and simple algorithm for training neural probabilistic language models, in: Proceedings of the 29th International Conference on International Conference on Machine Learning, ICML’12, Omnipress, USA, 2012, pp. 419–426. URL: <http://dl.acm.org/citation.cfm?id=3042573.3042630>.
 - [233] S. Jean, O. Firat, K. Cho, R. Memisevic, Y. Bengio, Montreal neural machine translation systems for WMT’15, in: Proceedings of the Tenth Workshop on Statistical Machine Translation, Association for Computational Linguistics, Lisbon, Portugal, 2015, pp. 134–140. URL: <https://www.aclweb.org/anthology/W15-3014>. doi:10.18653/v1/W15-3014.
 - [234] H. Mi, Z. Wang, A. Ittycheriah, Vocabulary manipulation for neural machine translation, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Association for Computational Linguistics, Berlin, Germany, 2016, pp. 124–129. URL: <https://www.aclweb.org/anthology/P16-2021>. doi:10.18653/v1/P16-2021.
 - [235] F. Stahlberg, E. Hasler, A. Waite, B. Byrne, Syntactically guided neural machine translation, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Association for Computational Linguistics, Berlin, Germany, 2016, pp. 299–305. URL: <https://www.aclweb.org/anthology/P16-2049>. doi:10.18653/v1/P16-2049.
 - [236] B. Sankaran, M. Freitag, Y. Al-Onaizan, Attention-based vocabulary selection for NMT decoding, arXiv preprint arXiv:1706.03824 (2017).
 - [237] G. L’Hostis, D. Grangier, M. Auli, Vocabulary selection strategies for neural machine translation, arXiv preprint arXiv:1610.00072 (2016).
 - [238] K. Hans, R. Milton, Improving the performance of neural machine translation involving morphologically rich languages, arXiv preprint arXiv:1612.02482 (2016).
 - [239] A. Tamchyna, M. Weller-Di Marco, A. Fraser, Modeling target-side inflection in neural machine translation, in: Proceedings of the Second Conference on Machine Translation, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 32–42. URL: <https://www.aclweb.org/anthology/W17-4704>. doi:10.18653/v1/W17-4704.
 - [240] W. Ling, I. Trancoso, C. Dyer, A. W. Black, Character-based neural machine translation, arXiv preprint arXiv:1511.04586 (2015).
 - [241] A. R. Johansen, J. M. Hansen, E. K. Obeid, C. K. Sønderby, O. Winther, Neural machine translation with characters and hierarchical encoding, arXiv preprint arXiv:1610.06550 (2016).

- [242] M. R. Costa-jussà, J. A. Fonollosa, Character-based neural machine translation, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Association for Computational Linguistics, Berlin, Germany, 2016, pp. 357–361. URL: <https://www.aclweb.org/anthology/P16-2058>. doi:10.18653/v1/P16-2058.
- [243] Y. Kim, Y. Jernite, D. Sontag, A. M. Rush, Character-aware neural language models, in: Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, AAAI’16, AAAI Press, 2016, pp. 2741–2749. URL: <http://dl.acm.org/citation.cfm?id=3016100.3016285>.
- [244] M.-T. Luong, C. D. Manning, Achieving open vocabulary neural machine translation with hybrid word-character models, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Berlin, Germany, 2016, pp. 1054–1063. URL: <https://www.aclweb.org/anthology/P16-1100>. doi:10.18653/v1/P16-1100.
- [245] M. Domingo, M. Garcia-Martinez, A. Helle, F. Casacuberta, How much does tokenization affect in neural machine translation?, arXiv preprint arXiv:1812.08621 (2018).
- [246] J. Chung, K. Cho, Y. Bengio, A character-level decoder without explicit segmentation for neural machine translation, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Berlin, Germany, 2016, pp. 1693–1703. URL: <https://www.aclweb.org/anthology/P16-1160>. doi:10.18653/v1/P16-1160.
- [247] J. Lee, K. Cho, T. Hofmann, Fully character-level neural machine translation without explicit segmentation, Transactions of the Association for Computational Linguistics 5 (2017) 365–378.
- [248] Z. Yang, W. Chen, F. Wang, B. Xu, A character-aware encoder for neural machine translation, in: Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers, The COLING 2016 Organizing Committee, Osaka, Japan, 2016, pp. 3063–3070. URL: <https://www.aclweb.org/anthology/C16-1288>.
- [249] M. R. Costa-jussà, C. Escolano, J. A. Fonollosa, Byte-based neural machine translation, in: Proceedings of the First Workshop on Subword and Character Level Models in NLP, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 154–158. URL: <https://www.aclweb.org/anthology/W17-4123>. doi:10.18653/v1/W17-4123.
- [250] C. Gulcehre, F. Dutil, A. Trischler, Y. Bengio, Plan, attend, generate: Character-level neural machine translation with planning, in: Proceedings of the 2nd Workshop on Representation Learning for NLP, Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 228–234. URL: <https://www.aclweb.org/anthology/W17-2627>. doi:10.18653/v1/W17-2627.
- [251] R. Chitnis, J. DeNero, Variable-length word encodings for neural translation models, in: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Lisbon, Portugal, 2015, pp. 2088–2093. URL: <https://www.aclweb.org/anthology/D15-1249>. doi:10.18653/v1/D15-1249.
- [252] M. Schuster, K. Nakajima, Japanese and korean voice search, in: 2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2012, pp. 5149–5152. doi:10.1109/ICASSP.2012.6289079.
- [253] P. Gage, A new algorithm for data compression, The C Users Journal 12 (1994) 23–38.
- [254] A. Kunchukuttan, P. Bhattacharyya, Learning variable length units for SMT between related languages via byte pair encoding, in: Proceedings of the First Workshop on Subword and Character Level Models in NLP, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 14–24. URL: <https://www.aclweb.org/anthology/W17-4102>. doi:10.18653/v1/W17-4102.
- [255] A. Kunchukuttan, P. Bhattacharyya, Faster decoding for subword level phrase-based SMT between related languages, in: Proceedings of the Third Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial3), The COLING 2016 Organizing Committee, Osaka, Japan, 2016, pp. 82–88. URL: <https://www.aclweb.org/anthology/W16-4811>.
- [256] N. F. Liu, J. May, M. Pust, K. Knight, Augmenting statistical machine translation with subword translation of out-of-vocabulary words, arXiv preprint arXiv:1808.05700 (2018).
- [257] Y. Wu, H. Zhao, Finding better subword segmentation for neural machine translation, in: M. Sun, T. Liu, X. Wang, Z. Liu, Y. Liu (Eds.), Chinese Computational Linguistics and Natural Language Processing Based on Naturally Annotated Big Data, Springer International Publishing, Cham, 2018, pp. 53–64.
- [258] E. Salesky, A. Runge, A. Coda, J. Niehues, G. Neubig, Optimizing segmentation granularity for neural machine translation, arXiv preprint arXiv:1810.08641 (2018).
- [259] T. Kudo, Subword regularization: Improving neural network translation models with multiple subword candidates, in: Proceedings of the 56th Annual Meeting of the Association for Computational

- Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 66–75. URL: <https://www.aclweb.org/anthology/P18-1007>.
- [260] T. Kudo, J. Richardson, SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 66–71. URL: <https://www.aclweb.org/anthology/D18-2012>.
 - [261] M. Huck, S. Riess, A. Fraser, Target-side word segmentation strategies for neural machine translation, in: Proceedings of the Second Conference on Machine Translation, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 56–67. URL: <https://www.aclweb.org/anthology/W17-4706>. doi:10.18653/v1/W17-4706.
 - [262] D. Macháček, J. Vidra, O. Bojar, Morphological and language-agnostic word segmentation for NMT, in: P. Sojka, A. Horák, I. Kopeček, K. Pala (Eds.), Text, Speech, and Dialogue, Springer International Publishing, Cham, 2018, pp. 277–284.
 - [263] D. Ataman, M. Negri, M. Turchi, M. Federico, Linguistically motivated vocabulary reduction for neural machine translation from Turkish to English, The Prague Bulletin of Mathematical Linguistics 108 (2017) 331–342.
 - [264] M. Pinnis, R. Krišlauks, D. Dekšne, T. Miks, Neural machine translation for morphologically rich languages with improved sub-word units and synthetic data, in: K. Ekšteins, V. Matoušek (Eds.), Text, Speech, and Dialogue, Springer International Publishing, Cham, 2017, pp. 237–245.
 - [265] R. Sennrich, How grammatical is character-level neural machine translation? Assessing MT quality with contrastive translation pairs, in: Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers, Association for Computational Linguistics, Valencia, Spain, 2017, pp. 376–382. URL: <https://www.aclweb.org/anthology/E17-2060>.
 - [266] J. Kreutzer, A. Sokolov, Optimally segmenting inputs for NMT shows preference for character-level processing, arXiv preprint arXiv:1810.01480 (2018).
 - [267] N. Durrani, F. Dalvi, H. Sajjad, Y. Belinkov, P. Nakov, What is in a translation unit? Comparing character and subword representations beyond translation, openreview.net (2018).
 - [268] Y. Belinkov, Y. Bisk, Synthetic and natural noise both break neural machine translation, arXiv preprint arXiv:1711.02173 (2017).
 - [269] C. Cherry, G. Foster, A. Bapna, O. Firat, W. Macherey, Revisiting character-based neural machine translation with capacity and compression, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 4295–4305. URL: <https://www.aclweb.org/anthology/D18-1461>.
 - [270] P. F. Brown, S. A. Della Pietra, V. J. Della Pietra, R. L. Mercer, The mathematics of statistical machine translation: Parameter estimation, Computational Linguistics 19 (1993) 263–311.
 - [271] C. Gulcehre, O. Firat, K. Xu, K. Cho, L. Barrault, H.-C. Lin, F. Bougares, H. Schwenk, Y. Bengio, On using monolingual corpora in neural machine translation, arXiv preprint arXiv:1503.03535 (2015).
 - [272] C. Gulcehre, O. Firat, K. Xu, K. Cho, Y. Bengio, On integrating a language model into neural machine translation, Computer Speech & Language 45 (2017) 137 – 148.
 - [273] F. Stahlberg, J. Cross, V. Stoyanov, Simple fusion: Return of the language model, in: Proceedings of the Third Conference on Machine Translation: Research Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 204–211. URL: <https://www.aclweb.org/anthology/W18-6321>.
 - [274] R. Sennrich, B. Haddow, A. Birch, Improving neural machine translation models with monolingual data, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Berlin, Germany, 2016, pp. 86–96. URL: <https://www.aclweb.org/anthology/P16-1009>. doi:10.18653/v1/P16-1009.
 - [275] A. Currey, A. V. Miceli Barone, K. Heafield, Copied monolingual data improves low-resource neural machine translation, in: Proceedings of the Second Conference on Machine Translation, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 148–156. URL: <https://www.aclweb.org/anthology/W17-4715>. doi:10.18653/v1/W17-4715.
 - [276] H. Schwenk, Investigations on large-scale lightly-supervised training for statistical machine translation., in: International Workshop on Spoken Language Translation (IWSLT) 2008, 2008, pp. 182–189.
 - [277] F. Burlot, F. Yvon, Using monolingual data in neural machine translation: A systematic study, in: Proceedings of the Third Conference on Machine Translation: Research Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 144–155. URL: <https://www.aclweb.org/anthology/W18-6321>.

- org/anthology/W18-6315.
- [278] V. C. D. Hoang, P. Koehn, G. Haffari, T. Cohn, Iterative back-translation for neural machine translation, in: *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 18–24. URL: <https://www.aclweb.org/anthology/W18-2703>.
 - [279] X. Niu, M. Denkowski, M. Carpuat, Bi-Directional neural machine translation with synthetic parallel data, in: *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 84–91. URL: <https://www.aclweb.org/anthology/W18-2710>.
 - [280] Z. Zhang, S. Liu, M. Li, M. Zhou, E. Chen, Joint training for neural machine translation models with monolingual data, in: *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
 - [281] A. Poncelas, D. Shterionov, A. Way, G. M. d. B. Wenniger, P. Passban, Investigating backtranslation in neural machine translation, arXiv preprint arXiv:1804.06189 (2018).
 - [282] S. Edunov, M. Ott, M. Auli, D. Grangier, Understanding back-translation at scale, in: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 489–500. URL: <https://www.aclweb.org/anthology/D18-1045>.
 - [283] X. Wang, H. Pham, Z. Dai, G. Neubig, SwitchOut: An efficient data augmentation algorithm for neural machine translation, in: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 856–861. URL: <https://www.aclweb.org/anthology/D18-1100>.
 - [284] K. Imamura, A. Fujita, E. Sumita, Enhancement of encoder and attention using target monolingual corpora in neural machine translation, in: *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 55–63. URL: <https://www.aclweb.org/anthology/W18-2707>.
 - [285] Y. Cheng, W. Xu, Z. He, W. He, H. Wu, M. Sun, Y. Liu, Semi-supervised learning for neural machine translation, in: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Association for Computational Linguistics, Berlin, Germany, 2016, pp. 1965–1974. URL: <https://www.aclweb.org/anthology/P16-1185>. doi:10.18653/v1/P16-1185.
 - [286] Z. Tu, Y. Liu, L. Shang, X. Liu, H. Li, Neural machine translation with reconstruction, in: *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, AAAI’17, AAAI Press, 2017, pp. 3097–3103. URL: <http://dl.acm.org/citation.cfm?id=3298483.3298684>.
 - [287] C. Escolano, M. R. Costa-jussà, J. A. Fonollosa, (self-attentive) autoencoder-based universal language representation for machine translation, arXiv preprint arXiv:1810.06351 (2018).
 - [288] D. He, Y. Xia, T. Qin, L. Wang, N. Yu, T.-Y. Liu, W.-Y. Ma, Dual learning for machine translation, in: D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, R. Garnett (Eds.), *Advances in Neural Information Processing Systems 29*, Curran Associates, Inc., 2016, pp. 820–828. URL: <http://papers.nips.cc/paper/6469-dual-learning-for-machine-translation.pdf>.
 - [289] H. Hassan, A. Aue, C. Chen, V. Chowdhary, J. H. Clark, C. Federmann, X. Huang, M. Junczys-Dowmunt, W. D. Lewis, M. Li, et al., Achieving human parity on automatic Chinese to English news translation, arXiv preprint arXiv:1803.05567 (2018).
 - [290] Y. Wang, Y. Xia, L. Zhao, J. Bian, T. Qin, G. Liu, T.-Y. Liu, Dual transfer learning for neural machine translation with marginal distribution regularization, in: *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
 - [291] J. Zhang, C. Zong, Exploiting source-side monolingual data in neural machine translation, in: *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Austin, Texas, 2016, pp. 1535–1545. URL: <https://www.aclweb.org/anthology/D16-1160>. doi:10.18653/v1/D16-1160.
 - [292] T. Domhan, F. Hieber, Using target-side monolingual data for neural machine translation through multi-task learning, in: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 1500–1505. URL: <https://www.aclweb.org/anthology/D17-1158>. doi:10.18653/v1/D17-1158.
 - [293] P. Ramachandran, P. J. Liu, Q. V. Le, Unsupervised pretraining for sequence to sequence learning, in: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 383–391. URL: <https://www.aclweb.org/anthology/D17-1039>. doi:10.18653/v1/D17-1039.
 - [294] I. Skorokhodov, A. Rykachevskiy, D. Emelyanenko, S. Slotin, A. Ponkratov, Semi-supervised neural machine translation with language models, in: *Proceedings of the AMTA 2018 Workshop on*

- Technologies for MT of Low Resource Languages (LoResMT 2018), Association for Machine Translation in the Americas, Boston, MA, 2018, pp. 37–44. URL: <https://www.aclweb.org/anthology/W18-2205>.
- [295] J. Niehues, E. Cho, T.-L. Ha, A. Waibel, Analyzing neural MT search and model performance, in: Proceedings of the First Workshop on Neural Machine Translation, Association for Computational Linguistics, Vancouver, 2017, pp. 11–17. URL: <https://www.aclweb.org/anthology/W17-3202>. doi:10.18653/v1/W17-3202.
 - [296] K. Papineni, S. Roukos, T. Ward, W.-J. Zhu, Bleu: A method for automatic evaluation of machine translation, in: Proceedings of 40th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, Philadelphia, Pennsylvania, USA, 2002, pp. 311–318. URL: <https://www.aclweb.org/anthology/P02-1040>. doi:10.3115/1073083.1073135.
 - [297] K. Murray, D. Chiang, Correcting length bias in neural machine translation, in: Proceedings of the Third Conference on Machine Translation: Research Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 212–223. URL: <https://www.aclweb.org/anthology/W18-6322>.
 - [298] A. Kumar, S. Sarawagi, Calibration of encoder decoder models for neural machine translation, arXiv preprint arXiv:1903.00802 (2019).
 - [299] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, Z. Wojna, Rethinking the inception architecture for computer vision, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 2818–2826.
 - [300] P. Koehn, R. Knowles, Six challenges for neural machine translation, in: Proceedings of the First Workshop on Neural Machine Translation, Association for Computational Linguistics, Vancouver, 2017, pp. 28–39. URL: <https://www.aclweb.org/anthology/W17-3204>. doi:10.18653/v1/W17-3204.
 - [301] N. Boulanger-Lewandowski, Y. Bengio, P. Vincent, Audio chord recognition with recurrent neural networks., in: ISMIR, Citeseer, 2013, pp. 335–340.
 - [302] F. J. Och, Minimum error rate training in statistical machine translation, in: Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, Sapporo, Japan, 2003, pp. 160–167. URL: <https://www.aclweb.org/anthology/P03-1021>. doi:10.3115/1075096.1075117.
 - [303] W. He, Z. He, H. Wu, H. Wang, Improved neural machine translation with SMT features, in: Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, AAAI’16, AAAI Press, 2016, pp. 151–157. URL: <http://dl.acm.org/citation.cfm?id=3015812.3015835>.
 - [304] L. Huang, K. Zhao, M. Ma, When to finish? Optimal beam search for neural text generation (modulo beam size), in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 2134–2139. URL: <https://www.aclweb.org/anthology/D17-1227>. doi:10.18653/v1/D17-1227.
 - [305] Y. Yang, L. Huang, M. Ma, Breaking the beam search curse: A study of (re-)scoring methods and stopping criteria for neural machine translation, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 3054–3059. URL: <https://www.aclweb.org/anthology/D18-1342>.
 - [306] Z. Tu, Z. Lu, Y. Liu, X. Liu, H. Li, Modeling coverage for neural machine translation, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Berlin, Germany, 2016, pp. 76–85. URL: <https://www.aclweb.org/anthology/P16-1008>. doi:10.18653/v1/P16-1008.
 - [307] J. Yang, B. Zhang, Y. Qin, X. Zhang, Q. Lin, J. Su, Otem&utem: Over- and under-translation evaluation metric for nmt, in: M. Zhang, V. Ng, D. Zhao, S. Li, H. Zan (Eds.), Natural Language Processing and Chinese Computing, Springer International Publishing, Cham, 2018, pp. 291–302.
 - [308] X. Kong, Z. Tu, S. Shi, E. Hovy, T. Zhang, Neural machine translation with adequacy-oriented learning, arXiv preprint arXiv:1811.08541 (2018).
 - [309] Y. Li, T. Xiao, Y. Li, Q. Wang, C. Xu, J. Zhu, A simple and effective approach to coverage-aware neural machine translation, in: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 292–297. URL: <https://www.aclweb.org/anthology/P18-2047>.
 - [310] H. Mi, B. Sankaran, Z. Wang, A. Ittycheriah, Coverage embedding models for neural machine translation, in: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Austin, Texas, 2016, pp. 955–960. URL: <https://www.aclweb.org/anthology/D16-1096>. doi:10.18653/v1/D16-1096.
 - [311] C. Malaviya, P. Ferreira, A. F. T. Martins, Sparse and constrained attention for neural machine

- translation, in: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 370–376. URL: <https://www.aclweb.org/anthology/P18-2059>.
- [312] M. B. Kazimi, M. R. Costa-Jussá, Coverage for character based neural machine translation, *Procesamiento del Lenguaje Natural* 59 (2017) 99–106.
 - [313] A. Fan, D. Grangier, M. Auli, Controllable abstractive summarization, in: Proceedings of the 2nd Workshop on Neural Machine Translation and Generation, Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 45–54. URL: <https://www.aclweb.org/anthology/W18-2706>.
 - [314] Y. Liu, Z. Luo, K. Zhu, Controlling length in abstractive summarization using a convolutional neural network, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 4110–4119. URL: <https://www.aclweb.org/anthology/D18-1444>.
 - [315] Y. Kikuchi, G. Neubig, R. Sasano, H. Takamura, M. Okumura, Controlling output length in neural encoder-decoders, in: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Austin, Texas, 2016, pp. 1328–1338. URL: <https://www.aclweb.org/anthology/D16-1140>. doi:10.18653/v1/D16-1140.
 - [316] S. Takase, N. Okazaki, Positional encoding to control output sequence length, 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), Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 3999–4004. URL: <https://www.aclweb.org/anthology/N19-1401>.
 - [317] M. D. Zeiler, ADADELTA: An adaptive learning rate method, arXiv preprint arXiv:1212.5701 (2012).
 - [318] M. Ranzato, S. Chopra, M. Auli, W. Zaremba, Sequence level training with recurrent neural networks, arXiv preprint arXiv:1511.06732 (2015).
 - [319] S. Edunov, M. Ott, M. Auli, D. Grangier, M. Ranzato, Classical structured prediction losses for sequence to sequence learning, in: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), Association for Computational Linguistics, New Orleans, Louisiana, 2018, pp. 355–364. URL: <https://www.aclweb.org/anthology/N18-1033>. doi:10.18653/v1/N18-1033.
 - [320] R. Pascanu, T. Mikolov, Y. Bengio, On the difficulty of training recurrent neural networks, in: S. Dasgupta, D. McAllester (Eds.), Proceedings of the 30th International Conference on Machine Learning, volume 28 of *Proceedings of Machine Learning Research*, PMLR, Atlanta, Georgia, USA, 2013, pp. 1310–1318. URL: <http://proceedings.mlr.press/v28/pascanu13.html>.
 - [321] R. Livni, S. Shalev-Shwartz, O. Shamir, On the computational efficiency of training neural networks, in: Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, K. Q. Weinberger (Eds.), Advances in Neural Information Processing Systems 27, Curran Associates, Inc., 2014, pp. 855–863. URL: <http://papers.nips.cc/paper/5267-on-the-computational-efficiency-of-training-neural-networks.pdf>.
 - [322] B. Byrne, Generalization and maximum likelihood from small data sets, in: Neural Networks for Signal Processing III - Proceedings of the 1993 IEEE-SP Workshop, 1993, pp. 197–206. doi:10.1109/NNSP.1993.471869.
 - [323] J. Chorowski, N. Jaitly, Towards better decoding and language model integration in sequence to sequence models, in: Proc. Interspeech 2017, 2017, pp. 523–527. URL: <http://dx.doi.org/10.21437/Interspeech.2017-343>. doi:10.21437/Interspeech.2017-343.
 - [324] G. Pereyra, G. Tucker, J. Chorowski, Ł. Kaiser, G. E. Hinton, Regularizing neural networks by penalizing confident output distributions, arXiv preprint arXiv:1701.06548 (2017).
 - [325] S. McCandlish, J. Kaplan, D. Amodei, O. D. Team, An empirical model of large-batch training, arXiv preprint arXiv:1812.06162 (2018).
 - [326] S. L. Smith, P.-J. Kindermans, C. Ying, Q. V. Le, Don’t decay the learning rate, increase the batch size, arXiv preprint arXiv:1711.00489 (2017).
 - [327] M. Neishi, J. Sakuma, S. Tohda, S. Ishiwatari, N. Yoshinaga, M. Toyoda, A bag of useful tricks for practical neural machine translation: Embedding layer initialization and large batch size, in: Proceedings of the 4th Workshop on Asian Translation (WAT2017), Asian Federation of Natural Language Processing, Taipei, Taiwan, 2017, pp. 99–109. URL: <https://www.aclweb.org/anthology/W17-5708>.
 - [328] M. Morishita, Y. Oda, G. Neubig, K. Yoshino, K. Sudoh, S. Nakamura, An empirical study of mini-batch creation strategies for neural machine translation, in: Proceedings of the First Workshop

- on Neural Machine Translation, Association for Computational Linguistics, Vancouver, 2017, pp. 61–68. URL: <https://www.aclweb.org/anthology/W17-3208>. doi:10.18653/v1/W17-3208.
- [329] D. Saunders, F. Stahlberg, A. de Gispert, B. Byrne, Multi-representation ensembles and delayed SGD updates improve syntax-based NMT, in: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 319–325. URL: <https://www.aclweb.org/anthology/P18-2051>.
 - [330] S. Bengio, O. Vinyals, N. Jaitly, N. Shazeer, Scheduled sampling for sequence prediction with recurrent neural networks, in: C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, R. Garnett (Eds.), Advances in Neural Information Processing Systems 28, Curran Associates, Inc., 2015, pp. 1171–1179. URL: <http://papers.nips.cc/paper/5956-scheduled-sampling-for-sequence-prediction-with-recurrent-neural-networks.pdf>.
 - [331] Y. Keneshloo, T. Shi, C. K. Reddy, N. Ramakrishnan, Deep reinforcement learning for sequence to sequence models, arXiv preprint arXiv:1805.09461 (2018).
 - [332] V. Zhukov, E. Golikov, M. Kretov, Differentiable lower bound for expected bleu score, arXiv preprint arXiv:1712.04708 (2017).
 - [333] L. Wu, F. Tian, T. Qin, J. Lai, T.-Y. Liu, A study of reinforcement learning for neural machine translation, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 3612–3621. URL: <https://www.aclweb.org/anthology/D18-1397>.
 - [334] B. Zoph, Q. V. Le, Neural architecture search with reinforcement learning, arXiv preprint arXiv:1611.01578 (2016).
 - [335] S. Shen, Y. Cheng, Z. He, W. He, H. Wu, M. Sun, Y. Liu, Minimum risk training for neural machine translation, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Berlin, Germany, 2016, pp. 1683–1692. URL: <https://www.aclweb.org/anthology/P16-1159>. doi:10.18653/v1/P16-1159.
 - [336] Y. Xia, T. Qin, W. Chen, J. Bian, N. Yu, T.-Y. Liu, Dual supervised learning, in: Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML’17, JMLR.org, 2017, pp. 3789–3798. URL: <http://dl.acm.org/citation.cfm?id=3305890.3306073>.
 - [337] Y. Cheng, S. Shen, Z. He, W. He, H. Wu, M. Sun, Y. Liu, Agreement-based joint training for bidirectional attention-based neural machine translation, in: Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI’16, AAAI Press, 2016, pp. 2761–2767. URL: <http://dl.acm.org/citation.cfm?id=3060832.3061007>.
 - [338] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial nets, in: Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, K. Q. Weinberger (Eds.), Advances in Neural Information Processing Systems 27, Curran Associates, Inc., 2014, pp. 2672–2680. URL: <http://papers.nips.cc/paper/5423-generative-adversarial-nets.pdf>.
 - [339] Z. Zhang, S. Liu, M. Li, M. Zhou, E. Chen, Bidirectional generative adversarial networks for neural machine translation, in: Proceedings of the 22nd Conference on Computational Natural Language Learning, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 190–199. URL: <https://www.aclweb.org/anthology/K18-1019>.
 - [340] Z. Yang, W. Chen, F. Wang, B. Xu, Improving neural machine translation with conditional sequence generative adversarial nets, in: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), Association for Computational Linguistics, New Orleans, Louisiana, 2018, pp. 1346–1355. URL: <https://www.aclweb.org/anthology/N18-1122>. doi:10.18653/v1/N18-1122.
 - [341] L. Wu, Y. Xia, L. Zhao, F. Tian, T. Qin, J. Lai, T.-Y. Liu, Adversarial neural machine translation, arXiv preprint arXiv:1704.06933 (2017).
 - [342] L. Yu, W. Zhang, J. Wang, Y. Yu, SeqGAN: Sequence generative adversarial nets with policy gradient, in: Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, AAAI’17, AAAI Press, 2017, pp. 2852–2858. URL: <http://dl.acm.org/citation.cfm?id=3298483.3298649>.
 - [343] J. Li, W. Monroe, T. Shi, S. Jean, A. Ritter, D. Jurafsky, Adversarial learning for neural dialogue generation, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 2157–2169. URL: <https://www.aclweb.org/anthology/D17-1230>. doi:10.18653/v1/D17-1230.
 - [344] A. M. Lamb, A. G. A. P. Goyal, Y. Zhang, S. Zhang, A. C. Courville, Y. Bengio, Professor forcing: A new algorithm for training recurrent networks, in: D. D. Lee, M. Sugiyama,

- U. V. Luxburg, I. Guyon, R. Garnett (Eds.), *Advances in Neural Information Processing Systems* 29, Curran Associates, Inc., 2016, pp. 4601–4609. URL: <http://papers.nips.cc/paper/6099-professor-forcing-a-new-algorithm-for-training-recurrent-networks.pdf>.
- [345] M. Caccia, L. Caccia, W. Fedus, H. Larochelle, J. Pineau, L. Charlin, Language GANs falling short, arXiv preprint arXiv:1811.02549 (2018).
- [346] W. E. Zhang, Q. Z. Sheng, A. A. F. Alhazmi, Generating textual adversarial examples for deep learning models: A survey, arXiv preprint arXiv:1901.06796 (2019).
- [347] P. Michel, X. Li, G. Neubig, J. Pino, On evaluation of adversarial perturbations for sequence-to-sequence models, 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)*, Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 3103–3114. URL: <https://www.aclweb.org/anthology/N19-1314>.
- [348] M. Ribeiro, S. Singh, C. Guestrin, “Why should I trust you?”: Explaining the predictions of any classifier, in: *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations*, Association for Computational Linguistics, San Diego, California, 2016, pp. 97–101. URL: <https://www.aclweb.org/anthology/N16-3020>. doi:10.18653/v1/N16-3020.
- [349] F. Doshi-Velez, B. Kim, Towards a rigorous science of interpretable machine learning, arXiv preprint arXiv:1702.08608 (2017).
- [350] Z. C. Lipton, The mythos of model interpretability, *Queue* 16 (2018) 30:31–30:57.
- [351] G. Montavon, W. Samek, K.-R. Müller, Methods for interpreting and understanding deep neural networks, *Digital Signal Processing* 73 (2018) 1–15.
- [352] A. Alishahi, G. Chrupala, T. Linzen, Analyzing and interpreting neural networks for NLP: A report on the first BlackboxNLP workshop, *Natural Language Engineering* (2019).
- [353] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, W. Samek, On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation, *PLOS ONE* 10 (2015) 1–46.
- [354] R. Schwarzenberg, D. Harbecke, V. Macketanz, E. Avramidis, S. Möller, Train, sort, explain: Learning to diagnose translation models, arXiv preprint arXiv:1903.12017 (2019).
- [355] D. Alvarez-Melis, T. Jaakkola, A causal framework for explaining the predictions of black-box sequence-to-sequence models, in: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 412–421. URL: <https://www.aclweb.org/anthology/D17-1042>. doi:10.18653/v1/D17-1042.
- [356] X. Ma, K. Li, P. Koehn, An analysis of source context dependency in neural machine translation, in: *21st Annual Conference of the European Association for Machine Translation*, 2018, p. 189.
- [357] S. Feng, E. Wallace, A. Grissom II, M. Iyyer, P. Rodriguez, J. Boyd-Graber, Pathologies of neural models make interpretations difficult, in: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 3719–3728. URL: <https://www.aclweb.org/anthology/D18-1407>.
- [358] A. Karpathy, J. Johnson, L. Fei-Fei, Visualizing and understanding recurrent networks, arXiv preprint arXiv:1506.02078 (2015).
- [359] J. Li, X. Chen, E. Hovy, D. Jurafsky, Visualizing and understanding neural models in NLP, in: *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Association for Computational Linguistics, San Diego, California, 2016, pp. 681–691. URL: <https://www.aclweb.org/anthology/N16-1082>. doi:10.18653/v1/N16-1082.
- [360] Y. Ding, Y. Liu, H. Luan, M. Sun, Visualizing and understanding neural machine translation, in: *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 1150–1159. URL: <https://www.aclweb.org/anthology/P17-1106>. doi:10.18653/v1/P17-1106.
- [361] D. Cashman, G. Patterson, A. Mosca, N. Watts, S. Robinson, R. Chang, RNNbow: Visualizing learning via backpropagation gradients in RNNs, *IEEE Computer Graphics and Applications* 38 (2018) 39–50.
- [362] Y. Belinkov, N. Durrani, F. Dalvi, H. Sajjad, J. Glass, What do neural machine translation models learn about morphology?, in: *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 861–872. URL: <https://www.aclweb.org/anthology/P17-1080>. doi:10.18653/v1/P17-1080.

- [363] A. Bau, Y. Belinkov, H. Sajjad, N. Durrani, F. Dalvi, J. Glass, Identifying and controlling important neurons in neural machine translation, arXiv preprint arXiv:1811.01157 (2018).
- [364] F. Dalvi, A. Nortonsmith, A. Bau, Y. Belinkov, H. Sajjad, N. Durrani, J. Glass, NeuroX: A toolkit for analyzing individual neurons in neural networks, arXiv preprint arXiv:1812.09359 (2018).
- [365] F. Dalvi, N. Durrani, H. Sajjad, Y. Belinkov, A. Bau, J. Glass, What is one grain of sand in the desert? Analyzing individual neurons in deep NLP models, in: Proceedings of the AAAI Conference on Artificial Intelligence (AAAI), 2019.
- [366] G. Tang, R. Sennrich, J. Nivre, An analysis of attention mechanisms: The case of word sense disambiguation in neural machine translation, in: Proceedings of the Third Conference on Machine Translation: Research Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 26–35. URL: <https://www.aclweb.org/anthology/W18-6304>.
- [367] H. Ghader, C. Monz, What does attention in neural machine translation pay attention to?, in: Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Asian Federation of Natural Language Processing, Taipei, Taiwan, 2017, pp. 30–39. URL: <https://www.aclweb.org/anthology/I17-1004>.
- [368] M. Ott, M. Auli, D. Grangier, M. Ranzato, Analyzing uncertainty in neural machine translation, in: J. Dy, A. Krause (Eds.), Proceedings of the 35th International Conference on Machine Learning, volume 80 of *Proceedings of Machine Learning Research*, PMLR, StockholmÅdssan, Stockholm Sweden, 2018, pp. 3956–3965. URL: <http://proceedings.mlr.press/v80/ott18a.html>.
- [369] A. de Gispert, G. Blackwood, G. Iglesias, B. Byrne, N-gram posterior probability confidence measures for statistical machine translation: An empirical study, *Machine Translation* 27 (2013) 85–114.
- [370] N. Bach, F. Huang, Y. Al-Onaizan, Goodness: A method for measuring machine translation confidence, in: Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, Portland, Oregon, USA, 2011, pp. 211–219. URL: <https://www.aclweb.org/anthology/P11-1022>.
- [371] N. Ueffing, H. Ney, Word-level confidence estimation for machine translation using phrase-based translation models, in: Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Vancouver, British Columbia, Canada, 2005, pp. 763–770. URL: <https://www.aclweb.org/anthology/H05-1096>.
- [372] M. Rikters, M. Fishel, Confidence through attention, arXiv preprint arXiv:1710.03743 (2017).
- [373] M. Rikters, Debugging neural machine translations, arXiv preprint arXiv:1808.02733 (2018).
- [374] L. Specia, F. Blain, V. Logacheva, R. Astudillo, A. F. T. Martins, Findings of the WMT 2018 shared task on quality estimation, in: Proceedings of the Third Conference on Machine Translation: Shared Task Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 689–709. URL: <https://www.aclweb.org/anthology/W18-6451>.
- [375] L. Specia, Exploiting objective annotations for measuring translation post-editing effort, in: Proceedings of the 15th Conference of the European Association for Machine Translation, 2011, pp. 73–80.
- [376] S. Jain, B. C. Wallace, Attention is not Explanation, 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), Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 3543–3556. URL: <https://www.aclweb.org/anthology/N19-1357>.
- [377] S. Serrano, N. A. Smith, Is attention interpretable?, in: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Florence, Italy, 2019.
- [378] G. Brunner, Y. Liu, D. Pascual, O. Richter, R. Wattenhofer, On the validity of self-attention as explanation in transformer models, arXiv preprint arXiv:1908.04211 (2019).
- [379] E. Hasler, A. de Gispert, G. Iglesias, B. Byrne, Neural machine translation decoding with terminology constraints, in: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), Association for Computational Linguistics, New Orleans, Louisiana, 2018, pp. 506–512. URL: <https://www.aclweb.org/anthology/N18-2081>. doi:10.18653/v1/N18-2081.
- [380] H. Mi, Z. Wang, A. Ittycheriah, Supervised attentions for neural machine translation, in: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Austin, Texas, 2016, pp. 2283–2288. URL: <https://www.aclweb.org/anthology/D16-1249>. doi:10.18653/v1/D16-1249.
- [381] W. Chen, E. Matusov, S. Khadivi, J.-T. Peter, Guided alignment training for topic-aware neural

- machine translation, AMTA 2016, Vol. (2016) 121.
- [382] L. Liu, M. Utiyama, A. Finch, E. Sumita, Neural machine translation with supervised attention, in: Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers, The COLING 2016 Organizing Committee, Osaka, Japan, 2016, pp. 3093–3102. URL: <https://www.aclweb.org/anthology/C16-1291>.
 - [383] T. Alkhouli, H. Ney, Biasing attention-based recurrent neural networks using external alignment information, in: Proceedings of the Second Conference on Machine Translation, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 108–117. URL: <https://www.aclweb.org/anthology/W17-4711>. doi:10.18653/v1/W17-4711.
 - [384] T. Cohn, C. D. V. Hoang, E. Vymolova, K. Yao, C. Dyer, G. Haffari, Incorporating structural alignment biases into an attentional neural translation model, in: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, San Diego, California, 2016, pp. 876–885. URL: <https://www.aclweb.org/anthology/N16-1102>. doi:10.18653/v1/N16-1102.
 - [385] D. Lawson, C.-C. Chiu, G. Tucker, C. Raffel, K. Swersky, N. Jaitly, Learning hard alignments with variational inference, in: 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2018, pp. 5799–5803. doi:10.1109/ICASSP.2018.8461977.
 - [386] R. Aharoni, Y. Goldberg, Morphological inflection generation with hard monotonic attention, in: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 2004–2015. URL: <https://www.aclweb.org/anthology/P17-1183>. doi:10.18653/v1/P17-1183.
 - [387] C. Raffel, M.-T. Luong, P. J. Liu, R. J. Weiss, D. Eck, Online and linear-time attention by enforcing monotonic alignments, in: Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML’17, JMLR.org, 2017, pp. 2837–2846. URL: <http://dl.acm.org/citation.cfm?id=3305890.3305974>.
 - [388] L. Yu, J. Buys, P. Blunsom, Online segment to segment neural transduction, in: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Austin, Texas, 2016, pp. 1307–1316. URL: <https://www.aclweb.org/anthology/D16-1138>. doi:10.18653/v1/D16-1138.
 - [389] L. Yu, P. Blunsom, C. Dyer, E. Grefenstette, T. Kocisky, The neural noisy channel, arXiv preprint arXiv:1611.02554 (2016).
 - [390] E. Choi, D. Hewlett, J. Uszkoreit, I. Polosukhin, A. Lacoste, J. Berant, Coarse-to-fine question answering for long documents, in: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 209–220. URL: <https://www.aclweb.org/anthology/P17-1020>. doi:10.18653/v1/P17-1020.
 - [391] T. Shen, T. Zhou, G. Long, J. Jiang, S. Wang, C. Zhang, Reinforced self-attention network: A hybrid of hard and soft attention for sequence modeling, in: Proceedings of the 27th International Joint Conference on Artificial Intelligence, IJCAI’18, AAAI Press, 2018, pp. 4345–4352. URL: <http://dl.acm.org/citation.cfm?id=3304222.3304374>.
 - [392] T. Alkhouli, G. Bretschner, J.-T. Peter, M. Hethnawi, A. Guta, H. Ney, Alignment-based neural machine translation, in: Proceedings of the First Conference on Machine Translation, Association for Computational Linguistics, Berlin, Germany, 2016, pp. 54–65. URL: <https://www.aclweb.org/anthology/W16-2206>. doi:10.18653/v1/W16-2206.
 - [393] T. Alkhouli, G. Bretschner, H. Ney, On the alignment problem in multi-head attention-based neural machine translation, in: Proceedings of the Third Conference on Machine Translation: Research Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 177–185. URL: <https://www.aclweb.org/anthology/W18-6318>.
 - [394] T. Zenkel, J. Wuebker, J. DeNero, Adding interpretable attention to neural translation models improves word alignment, arXiv preprint arXiv:1901.11359 (2019).
 - [395] K. Tran, A. Bisazza, C. Monz, The importance of being recurrent for modeling hierarchical structure, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 4731–4736. URL: <https://www.aclweb.org/anthology/D18-1503>.
 - [396] T. Domhan, How much attention do you need? A granular analysis of neural machine translation architectures, in: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 1799–1808. URL: <https://www.aclweb.org/anthology/P18-1167>.
 - [397] K. Song, T. Xu, F. Peng, J. Lu, Hybrid self-attention network for machine translation, arXiv

- preprint arXiv:1811.00253 (2018).
- [398] K. Ahmed, N. S. Keskar, R. Socher, Weighted Transformer network for machine translation, arXiv preprint arXiv:1711.02132 (2017).
 - [399] Q. Guo, X. Qiu, P. Liu, Y. Shao, X. Xue, Z. Zhang, Star-Transformer, arXiv preprint arXiv:1902.09113 (2019).
 - [400] J. R. Medina, J. Kalita, Parallel attention mechanisms in neural machine translation, in: 2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA), 2018, pp. 547–552. doi:10.1109/ICMLA.2018.00088.
 - [401] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, Language models are unsupervised multitask learners, <https://openai.com/blog/better-language-models/>, 2019.
 - [402] Z. Dai, Z. Yang, Y. Yang, W. W. Cohen, J. Carbonell, Q. V. Le, R. Salakhutdinov, Transformer-XL: Attentive language models beyond a fixed-length context, arXiv preprint arXiv:1901.02860 (2019).
 - [403] B. Krause, E. Kahembwe, I. Murray, S. Renals, Dynamic evaluation of Transformer language models, arXiv preprint arXiv:1904.08378 (2019).
 - [404] F. Meng, Z. Tu, Y. Cheng, H. Wu, J. Zhai, Y. Yang, D. Wang, Neural machine translation with key-value memory-augmented attention, in: Proceedings of the 27th International Joint Conference on Artificial Intelligence, IJCAI’18, AAAI Press, 2018, pp. 2574–2580. URL: <http://dl.acm.org/citation.cfm?id=3304889.3305018>.
 - [405] Z. Yang, Z. Hu, Y. Deng, C. Dyer, A. Smola, Neural machine translation with recurrent attention modeling, in: Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers, Association for Computational Linguistics, Valencia, Spain, 2017, pp. 383–387. URL: <https://www.aclweb.org/anthology/E17-2061>.
 - [406] S. Feng, S. Liu, M. Li, M. Zhou, Implicit distortion and fertility models for attention-based encoder-decoder NMT model, arXiv preprint arXiv:1601.03317 (2016).
 - [407] H. Choi, K. Cho, Y. Bengio, Fine-grained attention mechanism for neural machine translation, *Neurocomputing* 284 (2018) 171 – 176.
 - [408] M. Rikters, O. Bojar, Paying attention to multi-word expressions in neural machine translation, arXiv preprint arXiv:1710.06313 (2017).
 - [409] S. Ishiwatari, J. Yao, S. Liu, M. Li, M. Zhou, N. Yoshinaga, M. Kitsuregawa, W. Jia, Chunk-based decoder for neural machine translation, in: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 1901–1912. URL: <https://www.aclweb.org/anthology/P17-1174>. doi:10.18653/v1/P17-1174.
 - [410] J. Feng, L. Kong, P.-S. Huang, C. Wang, D. Huang, J. Mao, K. Qiao, D. Zhou, Neural phrase-to-phrase machine translation, arXiv preprint arXiv:1811.02172 (2018).
 - [411] P.-S. Huang, C. Wang, S. Huang, D. Zhou, L. Deng, Towards neural phrase-based machine translation, arXiv preprint arXiv:1706.05565 (2017).
 - [412] Y. Li, D. Xiong, M. Zhang, Neural machine translation with phrasal attention, in: Machine Translation: 13th China Workshop, CWMt 2017, Revised Selected Papers, Springer, 2017, pp. 1–8.
 - [413] A. Eriguchi, K. Hashimoto, Y. Tsuruoka, Incorporating source-side phrase structures into neural machine translation, *Computational Linguistics* 45 (2019) 267–292.
 - [414] J. Lin, X. Sun, X. Ren, M. Li, Q. Su, Learning when to concentrate or divert attention: Self-adaptive attention temperature for neural machine translation, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 2985–2990. URL: <https://www.aclweb.org/anthology/D18-1331>.
 - [415] B. Zhang, D. Xiong, J. Su, A gru-gated attention model for neural machine translation, arXiv preprint arXiv:1704.08430 (2017).
 - [416] H. T. Siegelmann, E. D. Sontag, On the computational power of neural nets, *Journal of Computer and System Sciences* 50 (1995) 132–150.
 - [417] A. Graves, G. Wayne, I. Danihelka, Neural turing machines, arXiv preprint arXiv:1410.5401 (2014).
 - [418] E. Grefenstette, K. M. Hermann, M. Suleyman, P. Blunsom, Learning to transduce with unbounded memory, in: C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, R. Garnett (Eds.), *Advances in Neural Information Processing Systems* 28, Curran Associates, Inc., 2015, pp. 1828–1836. URL: <http://papers.nips.cc/paper/5648-learning-to-transduce-with-unbounded-memory.pdf>.

- [419] R. J. Williams, D. Zipser, A learning algorithm for continually running fully recurrent neural networks, *Neural computation* 1 (1989) 270–280.
- [420] G.-Z. Sun, H.-H. Chen, C. L. Giles, Y.-C. Lee, D. Chen, Connectionist pushdown automata that learn context-free grammars, in: *Proceedings of the International Joint Conference on Neural Networks*, volume 1, Lawrence Earlbaum Hillsdale, NJ, 1990, pp. 577–580.
- [421] G.-Z. Sun, C. L. Giles, H.-H. Chen, Y.-C. Lee, *The Neural Network Pushdown Automation: Model, Stack and Learning Simulations*, Technical Report, University of Maryland at College Park, College Park, MD, USA, 1993.
- [422] A. Joulin, T. Mikolov, Inferring algorithmic patterns with stack-augmented recurrent nets, in: C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, R. Garnett (Eds.), *Advances in Neural Information Processing Systems 28*, Curran Associates, Inc., 2015, pp. 190–198. URL: <http://papers.nips.cc/paper/5857-inferring-algorithmic-patterns-with-stack-augmented-recurrent-nets.pdf>.
- [423] K. Kurach, M. Andrychowicz, I. Sutskever, Neural random-access machines, *arXiv preprint arXiv:1511.06392* (2015).
- [424] A. Graves, G. Wayne, M. Reynolds, T. Harley, I. Danihelka, A. Grabska-Barwińska, S. G. Colmenarejo, E. Grefenstette, T. Ramalho, J. Agapiou, et al., Hybrid computing using a neural network with dynamic external memory, *Nature* 538 (2016) 471.
- [425] S. Chandar, S. Ahn, H. Larochelle, P. Vincent, G. Tesauro, Y. Bengio, Hierarchical memory networks, *arXiv preprint arXiv:1605.07427* (2016).
- [426] M. Wang, Z. Lu, H. Li, Q. Liu, Memory-enhanced decoder for neural machine translation, in: *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Austin, Texas, 2016, pp. 278–286. URL: <https://www.aclweb.org/anthology/D16-1027>. doi:10.18653/v1/D16-1027.
- [427] Y. Feng, S. Zhang, A. Zhang, D. Wang, A. Abel, Memory-augmented neural machine translation, in: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 1390–1399. URL: <https://www.aclweb.org/anthology/D17-1146>. doi:10.18653/v1/D17-1146.
- [428] Y. Li, X. Liu, D. Liu, X. Zhang, J. Liu, Learning efficient lexically-constrained neural machine translation with external memory, *arXiv preprint arXiv:1901.11344* (2019).
- [429] H. Xiong, Z. He, X. Hu, H. Wu, Multi-channel encoder for neural machine translation, in: *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [430] T. Vu, B. Hu, T. Munkhdalai, H. Yu, Sentence simplification with memory-augmented neural networks, in: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, Association for Computational Linguistics, New Orleans, Louisiana, 2018, pp. 79–85. URL: <https://www.aclweb.org/anthology/N18-2013>. doi:10.18653/v1/N18-2013.
- [431] S. Pramanik, A. Hussain, Text normalization using memory augmented neural networks, *Speech Communication* 109 (2019) 15 – 23.
- [432] B. Zhang, D. Xiong, J. Su, H. Duan, M. Zhang, Variational neural machine translation, in: *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Austin, Texas, 2016, pp. 521–530. URL: <https://www.aclweb.org/anthology/D16-1050>. doi:10.18653/v1/D16-1050.
- [433] J. Su, S. Wu, D. Xiong, Y. Lu, X. Han, B. Zhang, Variational recurrent neural machine translation, in: *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [434] J. Bastings, W. Aziz, I. Titov, K. Sima'an, Modeling latent sentence structure in neural machine translation, *arXiv preprint arXiv:1901.06436* (2019).
- [435] H. Shah, D. Barber, Generative neural machine translation, in: S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, R. Garnett (Eds.), *Advances in Neural Information Processing Systems 31*, Curran Associates, Inc., 2018, pp. 1346–1355. URL: <http://papers.nips.cc/paper/7409-generative-neural-machine-translation.pdf>.
- [436] C. Wang, J. Zhang, H. Chen, Semi-autoregressive neural machine translation, in: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 479–488. URL: <https://www.aclweb.org/anthology/D18-1044>.
- [437] J. Gu, J. Bradbury, C. Xiong, V. O. Li, R. Socher, Non-autoregressive neural machine translation, *arXiv preprint arXiv:1711.02281* (2017).
- [438] J. Guo, X. Tan, D. He, T. Qin, L. Xu, T.-Y. Liu, Non-autoregressive neural machine translation with enhanced decoder input, *arXiv preprint arXiv:1812.09664* (2018).

- [439] Y. Wang, F. Tian, D. He, T. Qin, C. Zhai, T.-Y. Liu, Non-autoregressive machine translation with auxiliary regularization, arXiv preprint arXiv:1902.10245 (2019).
- [440] J. Libovický, J. Helcl, End-to-end non-autoregressive neural machine translation with connectionist temporal classification, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 3016–3021. URL: <https://www.aclweb.org/anthology/D18-1336>.
- [441] J. Lee, E. Mansimov, K. Cho, Deterministic non-autoregressive neural sequence modeling by iterative refinement, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 1173–1182. URL: <https://www.aclweb.org/anthology/D18-1149>.
- [442] N. Akoury, K. Krishna, M. Iyyer, Syntactically supervised Transformers for faster neural machine translation, in: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Florence, Italy, 2019.
- [443] P. Bahar, C. Brix, H. Ney, Towards two-dimensional sequence to sequence model in neural machine translation, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 3009–3015. URL: <https://www.aclweb.org/anthology/D18-1335>.
- [444] L. Kaiser, S. Bengio, Can active memory replace attention?, in: D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, R. Garnett (Eds.), Advances in Neural Information Processing Systems 29, Curran Associates, Inc., 2016, pp. 3781–3789. URL: <http://papers.nips.cc/paper/6295-can-active-memory-replace-attention.pdf>.
- [445] T. He, X. Tan, Y. Xia, D. He, T. Qin, Z. Chen, T.-Y. Liu, Layer-wise coordination between encoder and decoder for neural machine translation, in: S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, R. Garnett (Eds.), Advances in Neural Information Processing Systems 31, Curran Associates, Inc., 2018, pp. 7944–7954. URL: <http://papers.nips.cc/paper/8019-layer-wise-coordination-between-encoder-and-decoder-for-neural-machine-translation.pdf>.
- [446] P. Resnik, Mining the web for bilingual text, in: Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, College Park, Maryland, USA, 1999, pp. 527–534. URL: <https://www.aclweb.org/anthology/P99-1068>. doi:10.3115/1034678.1034757.
- [447] P. Resnik, N. A. Smith, The web as a parallel corpus, American Journal of Computational Linguistics 29 (2003) 349–380.
- [448] H. Khayrallah, P. Koehn, On the impact of various types of noise on neural machine translation, in: Proceedings of the 2nd Workshop on Neural Machine Translation and Generation, Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 74–83. URL: <https://www.aclweb.org/anthology/W18-2709>.
- [449] S. Rarrick, C. Quirk, W. D. Lewis, MT detection in web-scraped parallel corpora, Proceedings of the Machine Translation Summit (MT Summit XIII) (2011).
- [450] Y. Arase, M. Zhou, Machine translation detection from monolingual web-text, in: Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Sofia, Bulgaria, 2013, pp. 1597–1607. URL: <https://www.aclweb.org/anthology/P13-1157>.
- [451] P. Michel, G. Neubig, MTNT: A testbed for machine translation of noisy text, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 543–553. URL: <https://www.aclweb.org/anthology/D18-1050>.
- [452] Y. Cheng, Z. Tu, F. Meng, J. Zhai, Y. Liu, Towards robust neural machine translation, in: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 1756–1766. URL: <https://www.aclweb.org/anthology/P18-1163>.
- [453] N. Ruiz, M. A. D. Gangi, N. Bertoldi, M. Federico, Assessing the tolerance of neural machine translation systems against speech recognition errors, in: Proc. Interspeech 2017, 2017, pp. 2635–2639. URL: <http://dx.doi.org/10.21437/Interspeech.2017-1690>. doi:10.21437/Interspeech.2017-1690.
- [454] V. Karpukhin, O. Levy, J. Eisenstein, M. Ghazvininejad, Training on synthetic noise improves robustness to natural noise in machine translation, arXiv preprint arXiv:1902.01509 (2019).
- [455] V. Vaibhav, S. Singh, C. Stewart, G. Neubig, Improving robustness of machine translation with

- synthetic noise, 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), Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 1916–1920. URL: <https://www.aclweb.org/anthology/N19-1190>.
- [456] K. Taghipour, S. Khadivi, J. Xu, Parallel corpus refinement as an outlier detection algorithm, Proceedings of the 13th Machine Translation Summit (MT Summit XIII) (2011) 414–421.
 - [457] L. Cui, D. Zhang, S. Liu, M. Li, M. Zhou, Bilingual data cleaning for SMT using graph-based random walk, in: Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Association for Computational Linguistics, Sofia, Bulgaria, 2013, pp. 340–345. URL: <https://www.aclweb.org/anthology/P13-2061>.
 - [458] A. Axelrod, X. He, J. Gao, Domain adaptation via pseudo in-domain data selection, in: Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Edinburgh, Scotland, UK., 2011, pp. 355–362. URL: <https://www.aclweb.org/anthology/D11-1033>.
 - [459] G. Foster, C. Goutte, R. Kuhn, Discriminative instance weighting for domain adaptation in statistical machine translation, in: Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Cambridge, MA, 2010, pp. 451–459. URL: <https://www.aclweb.org/anthology/D10-1044>.
 - [460] M. van der Wees, A. Bisazza, C. Monz, Dynamic data selection for neural machine translation, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 1400–1410. URL: <https://www.aclweb.org/anthology/D17-1147>. doi:10.18653/v1/D17-1147.
 - [461] M. Carpuat, Y. Vyas, X. Niu, Detecting cross-lingual semantic divergence for neural machine translation, in: Proceedings of the First Workshop on Neural Machine Translation, Association for Computational Linguistics, Vancouver, 2017, pp. 69–79. URL: <https://www.aclweb.org/anthology/W17-3209>. doi:10.18653/v1/W17-3209.
 - [462] P. Koehn, H. Khayrallah, K. Heafield, M. L. Forcada, Findings of the WMT 2018 shared task on parallel corpus filtering, in: Proceedings of the Third Conference on Machine Translation: Shared Task Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 726–739. URL: <https://www.aclweb.org/anthology/W18-6453>.
 - [463] M. Junczys-Dowmunt, Dual conditional cross-entropy filtering of noisy parallel corpora, in: Proceedings of the Third Conference on Machine Translation: Shared Task Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 888–895. URL: <https://www.aclweb.org/anthology/W18-6478>.
 - [464] N. Rossenbach, J. Rosendahl, Y. Kim, M. Graça, A. Gokrani, H. Ney, The RWTH Aachen University filtering system for the WMT 2018 parallel corpus filtering task, in: Proceedings of the Third Conference on Machine Translation: Shared Task Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 946–954. URL: <https://www.aclweb.org/anthology/W18-6487>.
 - [465] X. Xu, S. Kuang, D. Xiong, Two effective approaches to data reduction for neural machine translation: Static and dynamic sentence selection, in: 2018 International Conference on Asian Language Processing (IALP), 2018, pp. 159–164. doi:10.1109/IALP.2018.8629243.
 - [466] D. Zhang, J. Kim, J. Crego, J. Senellart, Boosting neural machine translation, in: Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers), Asian Federation of Natural Language Processing, Taipei, Taiwan, 2017, pp. 271–276. URL: <https://www.aclweb.org/anthology/I17-2046>.
 - [467] Y. Bengio, J. Louradour, R. Collobert, J. Weston, Curriculum learning, in: Proceedings of the 26th Annual International Conference on Machine Learning, ICML '09, ACM, 2009, pp. 41–48. URL: <http://doi.acm.org/10.1145/1553374.1553380>. doi:10.1145/1553374.1553380.
 - [468] W. Wang, T. Watanabe, M. Hughes, T. Nakagawa, C. Chelba, Denoising neural machine translation training with trusted data and online data selection, in: Proceedings of the Third Conference on Machine Translation: Research Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 133–143. URL: <https://www.aclweb.org/anthology/W18-6314>.
 - [469] G. Kumar, G. Foster, C. Cherry, M. Krikun, Reinforcement learning based curriculum optimization for 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), Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 2054–2061. URL: <https://www.aclweb.org/anthology/N19-1208>.
 - [470] E. A. Platanios, O. Stretcu, G. Neubig, B. Poczos, T. Mitchell, Competence-based curriculum

- learning for 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), Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 1162–1172. URL: <https://www.aclweb.org/anthology/N19-1119>.
- [471] C. Chu, R. Wang, A survey of domain adaptation for neural machine translation, in: Proceedings of the 27th International Conference on Computational Linguistics, Association for Computational Linguistics, Santa Fe, New Mexico, USA, 2018, pp. 1304–1319. URL: <https://www.aclweb.org/anthology/C18-1111>.
- [472] C. Chu, R. Dabre, S. Kurohashi, A comprehensive empirical comparison of domain adaptation methods for neural machine translation, *Journal of Information Processing* 26 (2018) 529–538.
- [473] A. S. Hildebrand, M. Eck, S. Vogel, A. Waibel, Adaptation of the translation model for statistical machine translation based on information retrieval, in: Proceedings of EAMT, volume 2005, 2005, pp. 133–142.
- [474] R. Wang, M. Utiyama, A. Finch, L. Liu, K. Chen, E. Sumita, Sentence selection and weighting for neural machine translation domain adaptation, *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 26 (2018) 1727–1741.
- [475] R. Wang, A. Finch, M. Utiyama, E. Sumita, Sentence embedding for neural machine translation domain adaptation, in: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 560–566. URL: <https://www.aclweb.org/anthology/P17-2089>. doi:10.18653/v1/P17-2089.
- [476] R. Wang, M. Utiyama, L. Liu, K. Chen, E. Sumita, Instance weighting for neural machine translation domain adaptation, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 1482–1488. URL: <https://www.aclweb.org/anthology/D17-1155>. doi:10.18653/v1/D17-1155.
- [477] B. Chen, C. Cherry, G. Foster, S. Larkin, Cost weighting for neural machine translation domain adaptation, in: Proceedings of the First Workshop on Neural Machine Translation, Association for Computational Linguistics, Vancouver, 2017, pp. 40–46. URL: <https://www.aclweb.org/anthology/W17-3205>. doi:10.18653/v1/W17-3205.
- [478] C. Kobus, J. Crego, J. Senellart, Domain control for neural machine translation, in: Proceedings of the International Conference Recent Advances in Natural Language Processing, RANLP 2017, INCOMA Ltd., Varna, Bulgaria, 2017, pp. 372–378. URL: https://doi.org/10.26615/978-954-452-049-6_049. doi:10.26615/978-954-452-049-6_049.
- [479] S. Tars, M. Fishel, Multi-domain neural machine translation, in: Proceedings of the 21st Annual Conference of the European Association for Machine Translation: 28-30 May 2018, Universitat d’Alacant, Alacant, Spain, European Association for Machine Translation, 2018, pp. 259–268.
- [480] D. Britz, Q. V. Le, R. Pryzant, Effective domain mixing for neural machine translation, in: Proceedings of the Second Conference on Machine Translation, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 118–126. URL: <https://www.aclweb.org/anthology/W17-4712>. doi:10.18653/v1/W17-4712.
- [481] H. Sajjad, N. Durrani, F. Dalvi, Y. Belinkov, S. Vogel, Neural machine translation training in a multi-domain scenario, in: International Workshop on Spoken Language Translation, 2017.
- [482] H. Khayrallah, G. Kumar, K. Duh, M. Post, P. Koehn, Neural lattice search for domain adaptation in machine translation, in: Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers), Asian Federation of Natural Language Processing, Taipei, Taiwan, 2017, pp. 20–25. URL: <https://www.aclweb.org/anthology/I17-2004>.
- [483] M.-T. Luong, C. D. Manning, Stanford neural machine translation systems for spoken language domains, in: Proceedings of the International Workshop on Spoken Language Translation, 2015, pp. 76–79.
- [484] R. M. French, Catastrophic forgetting in connectionist networks, *Trends in Cognitive Sciences* 3 (1999) 128 – 135.
- [485] I. Goodfellow, M. Mirza, D. Xiao, A. Courville, Y. Bengio, An empirical investigation of catastrophic forgetting in gradient-based neural networks, arXiv preprint arXiv:1312.6211 (2013).
- [486] B. Thompson, H. Khayrallah, A. Anastasopoulos, A. D. McCarthy, K. Duh, R. Marvin, P. McNamee, J. Gwinnup, T. Anderson, P. Koehn, Freezing subnetworks to analyze domain adaptation in neural machine translation, in: Proceedings of the Third Conference on Machine Translation: Research Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 124–132. URL: <https://www.aclweb.org/anthology/W18-6313>.
- [487] P. Swietojanski, S. Renals, Learning hidden unit contributions for unsupervised speaker adaptation

- of neural network acoustic models, in: 2014 IEEE Spoken Language Technology Workshop (SLT), 2014, pp. 171–176. doi:10.1109/SLT.2014.7078569.
- [488] D. Vilar, Learning hidden unit contribution for adapting neural machine translation models, in: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), Association for Computational Linguistics, New Orleans, Louisiana, 2018, pp. 500–505. URL: <https://www.aclweb.org/anthology/N18-2080>. doi:10.18653/v1/N18-2080.
 - [489] H. Khayrallah, B. Thompson, K. Duh, P. Koehn, Regularized training objective for continued training for domain adaptation in neural machine translation, in: Proceedings of the 2nd Workshop on Neural Machine Translation and Generation, Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 36–44. URL: <https://www.aclweb.org/anthology/W18-2705>.
 - [490] P. Dakwale, C. Monz, Fine-tuning for neural machine translation with limited degradation across in-and out-of-domain data, Proceedings of the XVI Machine Translation Summit (2017) 117.
 - [491] A. V. Miceli Barone, B. Haddow, U. Germann, R. Sennrich, Regularization techniques for fine-tuning in neural machine translation, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 1489–1494. URL: <https://www.aclweb.org/anthology/D17-1156>. doi:10.18653/v1/D17-1156.
 - [492] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, R. Hadsell, Overcoming catastrophic forgetting in neural networks, Proceedings of the National Academy of Sciences 114 (2017) 3521–3526.
 - [493] B. Thompson, J. Gwinnup, H. Khayrallah, K. Duh, P. Koehn, Overcoming catastrophic forgetting during domain adaptation of 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), Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 2062–2068. URL: <https://www.aclweb.org/anthology/N19-1209>.
 - [494] D. Saunders, A. de Gispert, F. Stahlberg, B. Byrne, Domain adaptive inference for neural machine translation, in: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Association for Computational Linguistics, 2019.
 - [495] S. Ren, W. Chen, S. Liu, M. Li, M. Zhou, S. Ma, Triangular architecture for rare language translation, in: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 56–65. URL: <https://www.aclweb.org/anthology/P18-1006>.
 - [496] B. Zoph, D. Yuret, J. May, K. Knight, Transfer learning for low-resource neural machine translation, in: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Austin, Texas, 2016, pp. 1568–1575. URL: <https://www.aclweb.org/anthology/D16-1163>. doi:10.18653/v1/D16-1163.
 - [497] T. Q. Nguyen, D. Chiang, Transfer learning across low-resource, related languages for neural machine translation, in: Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers), Asian Federation of Natural Language Processing, Taipei, Taiwan, 2017, pp. 296–301. URL: <https://www.aclweb.org/anthology/I17-2050>.
 - [498] R. Murthy, A. Kunchukuttan, P. Bhattacharyya, Addressing word-order divergence in multilingual neural machine translation for extremely low resource languages, 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), Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 3868–3873. URL: <https://www.aclweb.org/anthology/N19-1387>.
 - [499] R. Dabre, T. Nakagawa, H. Kazawa, An empirical study of language relatedness for transfer learning in neural machine translation, in: Proceedings of the 31st Pacific Asia Conference on Language, Information and Computation, The National University (Philippines), 2017, pp. 282–286. URL: <https://www.aclweb.org/anthology/Y17-1038>.
 - [500] G. Neubig, J. Hu, Rapid adaptation of neural machine translation to new languages, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 875–880. URL: <https://www.aclweb.org/anthology/D18-1103>.
 - [501] A. Tong, L. Diduch, J. Fiscus, Y. Haghpanah, S. Huang, D. Joy, K. Peterson, I. Soboroff, Overview of the NIST 2016 LoReHLT evaluation, Machine Translation 32 (2018) 11–30.
 - [502] A. Conneau, G. Lample, M. Ranzato, L. Denoyer, H. Jégou, Word translation without parallel

- data, arXiv preprint arXiv:1710.04087 (2017).
- [503] M. Artetxe, G. Labaka, E. Agirre, Learning bilingual word embeddings with (almost) no bilingual data, in: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 451–462. URL: <https://www.aclweb.org/anthology/P17-1042>. doi:10.18653/v1/P17-1042.
 - [504] Y. Hoshen, L. Wolf, Non-adversarial unsupervised word translation, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 469–478. URL: <https://www.aclweb.org/anthology/D18-1043>.
 - [505] G. Lample, A. Conneau, L. Denoyer, M. Ranzato, Unsupervised machine translation using monolingual corpora only, arXiv preprint arXiv:1711.00043 (2017).
 - [506] M. Artetxe, G. Labaka, E. Agirre, K. Cho, Unsupervised neural machine translation, arXiv preprint arXiv:1710.11041 (2017).
 - [507] G. Lample, M. Ott, A. Conneau, L. Denoyer, M. Ranzato, Phrase-based & neural unsupervised machine translation, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 5039–5049. URL: <https://www.aclweb.org/anthology/D18-1549>.
 - [508] S. Ren, Z. Zhang, S. Liu, M. Zhou, S. Ma, Unsupervised neural machine translation with SMT as posterior regularization, arXiv preprint arXiv:1901.04112 (2019).
 - [509] J. Wu, X. Wang, W. Y. Wang, Extract and edit: An alternative to back-translation for unsupervised 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), Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 1173–1183. URL: <https://www.aclweb.org/anthology/N19-1120>.
 - [510] M. Johnson, M. Schuster, Q. V. Le, M. Krikun, Y. Wu, Z. Chen, N. Thorat, F. Viégas, M. Wattenberg, G. Corrado, M. Hughes, J. Dean, Google’s multilingual neural machine translation system: Enabling zero-shot translation, Transactions of the Association for Computational Linguistics 5 (2017) 339–351.
 - [511] T.-L. Ha, J. Niehues, A. Waibel, Toward multilingual neural machine translation with universal encoder and decoder, in: International Workshop on Spoken Language Translation IWSLT, 2016.
 - [512] T.-L. Ha, J. Niehues, A. Waibel, Effective strategies in zero-shot neural machine translation, in: International Workshop on Spoken Language Translation, 2017.
 - [513] M.-T. Luong, Q. V. Le, I. Sutskever, O. Vinyals, L. Kaiser, Multi-task sequence to sequence learning, arXiv preprint arXiv:1511.06114 (2015).
 - [514] O. Firat, K. Cho, Y. Bengio, 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: Human Language Technologies, Association for Computational Linguistics, San Diego, California, 2016, pp. 866–875. URL: <https://www.aclweb.org/anthology/N16-1101>. doi:10.18653/v1/N16-1101.
 - [515] O. Firat, K. Cho, B. Sankaran, F. T. Y. Vural, Y. Bengio, Multi-way, multilingual neural machine translation, Computer Speech & Language 45 (2017) 236 – 252.
 - [516] D. Dong, H. Wu, W. He, D. Yu, H. Wang, Multi-task learning for multiple language translation, in: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Association for Computational Linguistics, Beijing, China, 2015, pp. 1723–1732. URL: <https://www.aclweb.org/anthology/P15-1166>. doi:10.3115/v1/P15-1166.
 - [517] Y. Cheng, Q. Yang, Y. Liu, M. Sun, W. Xu, Joint training for pivot-based neural machine translation, in: Proceedings of the 26th International Joint Conference on Artificial Intelligence, IJCAI’17, AAAI Press, 2017, pp. 3974–3980. URL: <http://dl.acm.org/citation.cfm?id=3171837.3171841>.
 - [518] Y. Lu, P. Keung, F. Ladhak, V. Bhardwaj, S. Zhang, J. Sun, A neural interlingua for multilingual machine translation, in: Proceedings of the Third Conference on Machine Translation: Research Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 84–92. URL: <https://www.aclweb.org/anthology/W18-6309>.
 - [519] S. M. Lakew, M. Cettolo, M. Federico, A comparison of Transformer and recurrent neural networks on multilingual neural machine translation, in: Proceedings of the 27th International Conference on Computational Linguistics, Association for Computational Linguistics, Santa Fe, New Mexico, USA, 2018, pp. 641–652. URL: <https://www.aclweb.org/anthology/C18-1054>.
 - [520] R. Aharoni, M. Johnson, O. Firat, Massively multilingual neural machine translation, in: Proceedings of the 2019 Conference of the North American Chapter of the Association for Compu-

- tational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 3874–3884. URL: <https://www.aclweb.org/anthology/N19-1388>.
- [521] M. Cettolo, M. Federico, L. Bentivogli, N. Jan, S. Sebastian, S. Katsutho, Y. Koichiro, F. Christian, Overview of the IWSLT 2017 evaluation campaign, in: International Workshop on Spoken Language Translation, 2017, pp. 2–14.
 - [522] F. J. Och, H. Ney, Statistical multi-source translation, in: Proceedings of MT Summit, volume 8, 2001, pp. 253–258.
 - [523] G.-H. Choi, J.-H. Shin, Y.-K. Kim, Improving a multi-source neural machine translation model with corpus extension for low-resource languages, in: Proceedings of the 11th Language Resources and Evaluation Conference, European Language Resource Association, Miyazaki, Japan, 2018. URL: <https://www.aclweb.org/anthology/L18-1144>.
 - [524] Y. Nishimura, K. Sudoh, G. Neubig, S. Nakamura, Multi-source neural machine translation with data augmentation, arXiv preprint arXiv:1810.06826 (2018).
 - [525] K. Kann, R. Cotterell, H. Schütze, Neural multi-source morphological reinflection, in: Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers, Association for Computational Linguistics, Valencia, Spain, 2017, pp. 514–524. URL: <https://www.aclweb.org/anthology/E17-1049>.
 - [526] A. Currey, K. Heafield, Multi-source syntactic neural machine translation, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 2961–2966. URL: <https://www.aclweb.org/anthology/D18-1327>.
 - [527] R. Bawden, R. Sennrich, A. Birch, B. Haddow, Evaluating discourse phenomena in neural machine translation, in: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), Association for Computational Linguistics, New Orleans, Louisiana, 2018, pp. 1304–1313. URL: <https://www.aclweb.org/anthology/N18-1118>. doi:10.18653/v1/N18-1118.
 - [528] R. Dabre, C. Chu, A. Kunchukuttan, A survey of multilingual neural machine translation, arXiv preprint arXiv:1905.05395 (2019).
 - [529] D. R. So, C. Liang, Q. V. Le, The evolved Transformer, arXiv preprint arXiv:1901.11117 (2019).
 - [530] H. Hoang, T. Dwojak, R. Krislauks, D. Torregrosa, K. Heafield, Fast neural machine translation implementation, in: Proceedings of the 2nd Workshop on Neural Machine Translation and Generation, Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 116–121. URL: <https://www.aclweb.org/anthology/W18-2714>.
 - [531] M. Ott, S. Edunov, D. Grangier, M. Auli, Scaling neural machine translation, in: Proceedings of the Third Conference on Machine Translation: Research Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 1–9. URL: <https://www.aclweb.org/anthology/W18-6301>.
 - [532] J. Quinn, M. Ballesteros, Pieces of eight: 8-bit neural machine translation, in: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 3 (Industry Papers), Association for Computational Linguistics, New Orleans - Louisiana, 2018, pp. 114–120. URL: <https://www.aclweb.org/anthology/N18-3014>. doi:10.18653/v1/N18-3014.
 - [533] J. Devlin, Sharp models on dull hardware: Fast and accurate neural machine translation decoding on the CPU, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 2820–2825. URL: <https://www.aclweb.org/anthology/D17-1300>. doi:10.18653/v1/D17-1300.
 - [534] J. Wu, C. Leng, Y. Wang, Q. Hu, J. Cheng, Quantized convolutional neural networks for mobile devices, in: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 4820–4828. doi:10.1109/CVPR.2016.521.
 - [535] Y. LeCun, J. S. Denker, S. A. Solla, R. E. Howard, L. D. Jackel, Optimal brain damage, in: Advances in neural information processing systems, volume 2, 1989, pp. 598–605.
 - [536] M. G. Augasta, T. Kathirvalavakumar, Pruning algorithms of neural networks — a comparative study, Central European Journal of Computer Science 3 (2013) 105–115.
 - [537] B. Hassibi, D. G. Stork, et al., Second order derivatives for network pruning: Optimal brain surgeon, Advances in neural information processing systems (1993) 164–171.
 - [538] S. Han, J. Pool, J. Tran, W. Dally, Learning both weights and connections for efficient neural network, in: C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, R. Garnett (Eds.), Advances in Neural Information Processing Systems 28, Cur-

- ran Associates, Inc., 2015, pp. 1135–1143. URL: <http://papers.nips.cc/paper/5784-learning-both-weights-and-connections-for-efficient-neural-network.pdf>.
- [539] A. See, M.-T. Luong, C. D. Manning, Compression of neural machine translation models via pruning, in: Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning, Association for Computational Linguistics, Berlin, Germany, 2016, pp. 291–301. URL: <https://www.aclweb.org/anthology/K16-1029>. doi:10.18653/v1/K16-1029.
 - [540] M. Zhu, S. Gupta, To prune, or not to prune: exploring the efficacy of pruning for model compression, arXiv preprint arXiv:1710.01878 (2017).
 - [541] S. Srinivas, R. V. Babu, Data-free parameter pruning for deep neural networks, in: Proceedings of the British Machine Vision Conference (BMVC), BMVA Press, 2015, pp. 31.1–31.12. URL: <https://dx.doi.org/10.5244/C.29.31>.
 - [542] M. Babaeizadeh, P. Smaragdis, R. H. Campbell, NoiseOut: A simple way to prune neural networks, in: Proceedings of the 1st International Workshop on Efficient Methods for Deep Neural Networks (EMDNN), 2016.
 - [543] M. Denil, B. Shakibi, L. Dinh, M. Ranzato, N. de Freitas, Predicting parameters in deep learning, in: C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, K. Q. Weinberger (Eds.), Advances in Neural Information Processing Systems 26, Curran Associates, Inc., 2013, pp. 2148–2156. URL: <http://papers.nips.cc/paper/5025-predicting-parameters-in-deep-learning.pdf>.
 - [544] E. L. Denton, W. Zaremba, J. Bruna, Y. LeCun, R. Fergus, Exploiting linear structure within convolutional networks for efficient evaluation, in: Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, K. Q. Weinberger (Eds.), Advances in Neural Information Processing Systems 27, Curran Associates, Inc., 2014, pp. 1269–1277. URL: <http://papers.nips.cc/paper/5544-exploiting-linear-structure-within-convolutional-networks-for-efficient-evaluation.pdf>.
 - [545] J. Xue, J. Li, Y. Gong, Restructuring of deep neural network acoustic models with singular value decomposition, in: Interspeech, 2013, pp. 2365–2369.
 - [546] R. Prabhavalkar, O. Alsharif, A. Bruguier, L. McGraw, On the compression of recurrent neural networks with an application to LVCSR acoustic modeling for embedded speech recognition, in: 2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2016, pp. 5970–5974. doi:10.1109/ICASSP.2016.7472823.
 - [547] Z. Lu, V. Sindhwani, T. N. Sainath, Learning compact recurrent neural networks, in: 2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2016, pp. 5960–5964. doi:10.1109/ICASSP.2016.7472821.
 - [548] C. Buciluă, R. Caruana, A. Niculescu-Mizil, Model compression, in: Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '06, ACM, New York, NY, USA, 2006, pp. 535–541. URL: <http://doi.acm.org/10.1145/1150402.1150464>. doi:10.1145/1150402.1150464.
 - [549] G. E. Hinton, O. Vinyals, J. Dean, Distilling the knowledge in a neural network, in: NIPS Deep Learning and Representation Learning Workshop, 2015. URL: <http://arxiv.org/abs/1503.02531>.
 - [550] J. H. Wong, M. J. Gales, Sequence student-teacher training of deep neural networks, in: Interspeech 2016, 2016, pp. 2761–2765. URL: <http://dx.doi.org/10.21437/Interspeech.2016-911>. doi:10.21437/Interspeech.2016-911.
 - [551] M. Freitag, Y. Al-Onaizan, B. Sankaran, Ensemble distillation for neural machine translation, arXiv preprint arXiv:1702.01802 (2017).
 - [552] D. Zhang, J. Crego, J. Senellart, Analyzing knowledge distillation in neural machine translation, in: International Workshop on Spoken Language Translation IWSLT, 2018.
 - [553] H.-G. Kim, H. Na, H. Lee, J. Lee, T. G. Kang, M.-J. Lee, Y. S. Choi, Knowledge distillation using output errors for self-attention end-to-end models, in: ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2019, pp. 6181–6185. doi:10.1109/ICASSP.2019.8682775.
 - [554] Y. Liu, H. Xiong, Z. He, J. Zhang, H. Wu, H. Wang, C. Zong, End-to-end speech translation with knowledge distillation, arXiv preprint arXiv:1904.08075 (2019).
 - [555] D. Elliott, S. Frank, E. Hasler, Multilingual image description with neural sequence models, arXiv preprint arXiv:1510.04709 (2015).
 - [556] J. Hitschler, S. Schamoni, S. Riezler, Multimodal pivots for image caption translation, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Berlin, Germany, 2016, pp. 2399–2409. URL: <https://www.aclweb.org/anthology/P16-1227>. doi:10.18653/v1/P16-1227.
 - [557] P.-Y. Huang, F. Liu, S.-R. Shiang, J. Oh, C. Dyer, Attention-based multimodal neural machine

- translation, in: Proceedings of the First Conference on Machine Translation, Association for Computational Linguistics, Berlin, Germany, 2016, pp. 639–645. URL: <https://www.aclweb.org/anthology/W16-2360>. doi:10.18653/v1/W16-2360.
- [558] L. Specia, S. Frank, K. Sima'an, D. Elliott, A shared task on multimodal machine translation and crosslingual image description, in: Proceedings of the First Conference on Machine Translation, Association for Computational Linguistics, Berlin, Germany, 2016, pp. 543–553. URL: <https://www.aclweb.org/anthology/W16-2346>. doi:10.18653/v1/W16-2346.
- [559] D. Elliott, S. Frank, L. Barrault, F. Bougares, L. Specia, Findings of the second shared task on multimodal machine translation and multilingual image description, in: Proceedings of the Second Conference on Machine Translation, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 215–233. URL: <https://www.aclweb.org/anthology/W17-4718>. doi:10.18653/v1/W17-4718.
- [560] L. Barrault, F. Bougares, L. Specia, C. Lala, D. Elliott, S. Frank, Findings of the third shared task on multimodal machine translation, in: Proceedings of the Third Conference on Machine Translation: Shared Task Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 304–323. URL: <https://www.aclweb.org/anthology/W18-6402>.
- [561] I. Calixto, Q. Liu, An error analysis for image-based multi-modal neural machine translation, Machine Translation (2019).
- [562] O. Vinyals, L. Kaiser, T. Koo, S. Petrov, I. Sutskever, G. E. Hinton, Grammar as a foreign language, in: C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, R. Garnett (Eds.), Advances in Neural Information Processing Systems 28, Curran Associates, Inc., 2015, pp. 2773–2781. URL: <http://papers.nips.cc/paper/5635-grammar-as-a-foreign-language.pdf>.
- [563] R. Aharoni, Y. Goldberg, Towards string-to-tree neural machine translation, in: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 132–140. URL: <https://www.aclweb.org/anthology/P17-2021>. doi:10.18653/v1/P17-2021.
- [564] C. Ma, L. Liu, A. Tamura, T. Zhao, E. Sumita, Deterministic attention for sequence-to-sequence constituent parsing, in: Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, AAAI'17, AAAI Press, 2017, pp. 3237–3243. URL: <http://dl.acm.org/citation.cfm?id=3298023.3298039>.
- [565] M. Nadejde, S. Reddy, R. Sennrich, T. Dwojak, M. Junczys-Dowmunt, P. Koehn, A. Birch, Predicting target language CCG supertags improves neural machine translation, in: Proceedings of the Second Conference on Machine Translation, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 68–79. URL: <https://www.aclweb.org/anthology/W17-4707>. doi:10.18653/v1/W17-4707.
- [566] C. Ma, A. Tamura, M. Utiyama, T. Zhao, E. Sumita, Forest-based neural machine translation, in: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 1253–1263. URL: <https://www.aclweb.org/anthology/P18-1116>.
- [567] P. Zaremoondi, G. Haffari, Incorporating syntactic uncertainty in neural machine translation with a forest-to-sequence model, in: Proceedings of the 27th International Conference on Computational Linguistics, Association for Computational Linguistics, Santa Fe, New Mexico, USA, 2018, pp. 1421–1429. URL: <https://www.aclweb.org/anthology/C18-1120>.
- [568] C. Dyer, A. Kuncoro, M. Ballesteros, N. A. Smith, Recurrent neural network grammars, in: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, San Diego, California, 2016, pp. 199–209. URL: <https://www.aclweb.org/anthology/N16-1024>. doi:10.18653/v1/N16-1024.
- [569] A. Eriguchi, Y. Tsuruoka, K. Cho, Learning to parse and translate improves neural machine translation, in: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 72–78. URL: <https://www.aclweb.org/anthology/P17-2012>. doi:10.18653/v1/P17-2012.
- [570] J. Bradbury, R. Socher, Towards neural machine translation with latent tree attention, in: Proceedings of the 2nd Workshop on Structured Prediction for Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 12–16. URL: <https://www.aclweb.org/anthology/W17-4303>. doi:10.18653/v1/W17-4303.
- [571] X. Wang, H. Pham, P. Yin, G. Neubig, A tree-based decoder for neural machine translation, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing,

- Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 4772–4777. URL: <https://www.aclweb.org/anthology/D18-1509>.
- [572] S. Wu, D. Zhang, N. Yang, M. Li, M. Zhou, Sequence-to-dependency neural machine translation, in: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 698–707. URL: <https://www.aclweb.org/anthology/P17-1065>. doi:10.18653/v1/P17-1065.
 - [573] K. S. Tai, R. Socher, C. D. Manning, Improved semantic representations from tree-structured long short-term memory networks, in: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Association for Computational Linguistics, Beijing, China, 2015, pp. 1556–1566. URL: <https://www.aclweb.org/anthology/P15-1150>. doi:10.3115/v1/P15-1150.
 - [574] A. Eriguchi, K. Hashimoto, Y. Tsuruoka, Tree-to-sequence attentional neural machine translation, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Berlin, Germany, 2016, pp. 823–833. URL: <https://www.aclweb.org/anthology/P16-1078>. doi:10.18653/v1/P16-1078.
 - [575] B. Yang, D. F. Wong, T. Xiao, L. S. Chao, J. Zhu, Towards bidirectional hierarchical representations for attention-based neural machine translation, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 1432–1441. URL: <https://www.aclweb.org/anthology/D17-1150>. doi:10.18653/v1/D17-1150.
 - [576] H. Chen, S. Huang, D. Chiang, J. Chen, Improved neural machine translation with a syntax-aware encoder and decoder, in: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 1936–1945. URL: <https://www.aclweb.org/anthology/P17-1177>. doi:10.18653/v1/P17-1177.
 - [577] Y. Shen, S. Tan, A. Sordoni, A. Courville, Ordered neurons: Integrating tree structures into recurrent neural networks, Proceedings of ICLR (2019).
 - [578] J. Bastings, I. Titov, W. Aziz, D. Marcheggiani, K. Simaan, Graph convolutional encoders for syntax-aware neural machine translation, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 1957–1967. URL: <https://www.aclweb.org/anthology/D17-1209>. doi:10.18653/v1/D17-1209.
 - [579] K. Chen, R. Wang, M. Utiyama, L. Liu, A. Tamura, E. Sumita, T. Zhao, Neural machine translation with source dependency representation, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 2846–2852. URL: <https://www.aclweb.org/anthology/D17-1304>. doi:10.18653/v1/D17-1304.
 - [580] K. Chen, R. Wang, M. Utiyama, E. Sumita, T. Zhao, Syntax-directed attention for neural machine translation, in: Thirty-Second AAAI Conference on Artificial Intelligence, 2018.
 - [581] Y. Kim, A. M. Rush, L. Yu, A. Kuncoro, C. Dyer, G. Melis, Unsupervised recurrent neural network grammars, arXiv preprint arXiv:1904.03746 (2019).
 - [582] J. Maillard, S. Clark, D. Yogatama, Jointly learning sentence embeddings and syntax with unsupervised Tree-LSTMs, arXiv preprint arXiv:1705.09189 (2017).
 - [583] A. Williams, A. Drozdov, S. Bowman, Do latent tree learning models identify meaningful structure in sentences?, Transactions of the Association for Computational Linguistics 6 (2018) 253–267.
 - [584] M. Sperber, G. Neubig, J. Niehues, A. Waibel, Neural lattice-to-sequence models for uncertain inputs, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 1380–1389. URL: <https://www.aclweb.org/anthology/D17-1145>. doi:10.18653/v1/D17-1145.
 - [585] J. Su, Z. Tan, D. Xiong, R. Ji, X. Shi, Y. Liu, Lattice-based recurrent neural network encoders for neural machine translation, in: Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, AAAI’17, AAAI Press, 2017, pp. 3302–3308. URL: <http://dl.acm.org/citation.cfm?id=3298023.3298048>.
 - [586] Z. Tan, J. Su, B. Wang, Y. Chen, X. Shi, Lattice-to-sequence attentional neural machine translation models, Neurocomputing 284 (2018) 138 – 147.
 - [587] D. Moussallem, M. Arčan, A.-C. N. Ngomo, P. Buitelaar, Augmenting neural machine translation with knowledge graphs, arXiv preprint arXiv:1902.08816 (2019).
 - [588] R. Koncel-Kedziorski, D. Bekal, Y. Luan, M. Lapata, H. Hajishirzi, Text Generation from Knowledge Graphs with Graph Transformers, 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), Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 2284–2293. URL: <https://www.aclweb.org/anthology/N19-1238>.
- [589] D. Marcheggiani, J. Bastings, I. Titov, Exploiting semantics in neural machine translation with graph convolutional networks, in: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), Association for Computational Linguistics, New Orleans, Louisiana, 2018, pp. 486–492. URL: <https://www.aclweb.org/anthology/N18-2078>. doi:10.18653/v1/N18-2078.
 - [590] P. Koehn, H. Hoang, Factored translation models, in: Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL), Association for Computational Linguistics, Prague, Czech Republic, 2007, pp. 868–876. URL: <https://www.aclweb.org/anthology/D07-1091>.
 - [591] R. Sennrich, B. Haddow, Linguistic input features improve neural machine translation, in: Proceedings of the First Conference on Machine Translation, Association for Computational Linguistics, Berlin, Germany, 2016, pp. 83–91. URL: <https://www.aclweb.org/anthology/W16-2209>. doi:10.18653/v1/W16-2209.
 - [592] M. García-Martínez, L. Barrault, F. Bougares, Factored neural machine translation architectures, in: International Workshop on Spoken Language Translation (IWSLT’16), 2016.
 - [593] M. García-Martínez, L. Barrault, F. Bougares, Neural machine translation by generating multiple linguistic factors, in: N. Camelin, Y. Estève, C. Martín-Vide (Eds.), Statistical Language and Speech Processing, Springer International Publishing, Cham, 2017, pp. 21–31.
 - [594] S. Läubli, R. Sennrich, M. Volk, Has machine translation achieved human parity? a case for document-level evaluation, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 4791–4796. URL: <https://www.aclweb.org/anthology/D18-1512>.
 - [595] L. Wang, Z. Tu, A. Way, Q. Liu, Exploiting cross-sentence context for neural machine translation, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 2826–2831. URL: <https://www.aclweb.org/anthology/D17-1301>. doi:10.18653/v1/D17-1301.
 - [596] S. Jean, S. Lauly, O. Firat, K. Cho, Does neural machine translation benefit from larger context?, arXiv preprint arXiv:1704.05135 (2017).
 - [597] Z. Tu, Y. Liu, S. Shi, T. Zhang, Learning to remember translation history with a continuous cache, Transactions of the Association for Computational Linguistics 6 (2018) 407–420.
 - [598] S. Maruf, G. Haffari, Document context neural machine translation with memory networks, in: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 1275–1284. URL: <https://www.aclweb.org/anthology/P18-1118>.
 - [599] S. Kuang, D. Xiong, W. Luo, G. Zhou, Cache-based document-level neural machine translation, arXiv preprint arXiv:1711.11221 (2017).
 - [600] F. Stahlberg, D. Saunders, A. de Gispert, B. Byrne, CUED@WMT19:EWC&LMs, in: Proceedings of the Fourth Conference on Machine Translation: Shared Task Papers, Association for Computational Linguistics, 2019.
 - [601] L. Miculicich, D. Ram, N. Pappas, J. Henderson, Document-level neural machine translation with hierarchical attention networks, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 2947–2954. URL: <https://www.aclweb.org/anthology/D18-1325>.
 - [602] S. Maruf, A. F. T. Martins, G. Haffari, Selective attention for context-aware 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), Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 3092–3102. URL: <https://www.aclweb.org/anthology/N19-1313>.
 - [603] H. Xiong, Z. He, H. Wu, H. Wang, Modeling coherence for discourse neural machine translation, arXiv preprint arXiv:1811.05683 (2018).
 - [604] J. Tiedemann, Y. Scherrer, Neural machine translation with extended context, in: Proceedings of the Third Workshop on Discourse in Machine Translation, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 82–92. URL: <https://www.aclweb.org/anthology/W17-4811>. doi:10.18653/v1/W17-4811.
 - [605] E. Voita, P. Serdyukov, R. Sennrich, I. Titov, Context-aware neural machine translation learns anaphora resolution, in: Proceedings of the 56th Annual Meeting of the Association for Com-

- putational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 1264–1274. URL: <https://www.aclweb.org/anthology/P18-1117>.
- [606] J. Zhang, H. Luan, M. Sun, F. Zhai, J. Xu, M. Zhang, Y. Liu, Improving the Transformer translation model with document-level context, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 533–542. URL: <https://www.aclweb.org/anthology/D18-1049>.
 - [607] M. Müller, A. Rios, E. Voita, R. Sennrich, A large-scale test set for the evaluation of context-aware pronoun translation in neural machine translation, in: Proceedings of the Third Conference on Machine Translation: Research Papers, Association for Computational Linguistics, Belgium, Brussels, 2018, pp. 61–72. URL: <https://www.aclweb.org/anthology/W18-6307>.
 - [608] A. Toral, V. M. Sánchez-Cartagena, A multifaceted evaluation of neural versus phrase-based machine translation for 9 language directions, in: Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers, Association for Computational Linguistics, Valencia, Spain, 2017, pp. 1063–1073. URL: <https://www.aclweb.org/anthology/E17-1100>.
 - [609] L. Bentivogli, A. Bisazza, M. Cettolo, M. Federico, Neural versus phrase-based machine translation quality: A case study, in: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Austin, Texas, 2016, pp. 257–267. URL: <https://www.aclweb.org/anthology/D16-1025>. doi:10.18653/v1/D16-1025.
 - [610] L. Bentivogli, A. Bisazza, M. Cettolo, M. Federico, Neural versus phrase-based mt quality: An in-depth analysis on English→German and English→French, *Computer Speech & Language* 49 (2018) 52 – 70.
 - [611] S. Castilho, J. Moorkens, F. Gaspari, R. Sennrich, V. Sosoni, P. Georgakopoulou, P. Lohar, A. Way, A. V. M. Barone, M. Gialama, A comparative quality evaluation of PBSMT and NMT using professional translators, Proceedings of Machine Translation Summit XVI, Nagoya, Japan (2017).
 - [612] L. Volkart, P. Bouillon, S. Girletti, Statistical vs. neural machine translation: A comparison of mth and deepl at swiss post’s language service, in: Proceedings of the 40th Conference Translating and the Computer, London, United-Kingdom, 2018, pp. 145–150. URL: <https://archive-ouverte.unige.ch/unige:111777>.
 - [613] S. K. Mahata, S. Mandal, D. Das, S. Bandyopadhyay, SMT vs NMT: a comparison over Hindi & Bengali simple sentences, arXiv preprint arXiv:1812.04898 (2018).
 - [614] S. Castilho, J. Moorkens, F. Gaspari, I. Calixto, J. Tinsley, A. Way, Is neural machine translation the new state of the art?, *The Prague Bulletin of Mathematical Linguistics* 108 (2017) 109–120.
 - [615] P. Isabelle, C. Cherry, G. Foster, A challenge set approach to evaluating machine translation, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 2486–2496. URL: <https://www.aclweb.org/anthology/D17-1263>. doi:10.18653/v1/D17-1263.
 - [616] C. Schnober, S. Eger, E.-L. Do Dinh, I. Gurevych, Still not there? Comparing traditional sequence-to-sequence models to encoder-decoder neural networks on monotone string translation tasks, in: Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers, The COLING 2016 Organizing Committee, Osaka, Japan, 2016, pp. 1703–1714. URL: <https://www.aclweb.org/anthology/C16-1160>.
 - [617] M. A. Menacer, D. Langlois, O. Mella, D. Fohr, D. Juvet, K. Smaïli, Is statistical machine translation approach dead?, in: ICNLSSP 2017 - International Conference on Natural Language, Signal and Speech Processing, ISGA, Casablanca, Morocco, 2017, pp. 1–5. URL: <https://hal.inria.fr/hal-01660016>.
 - [618] M. Dowling, T. Lynn, A. Poncelas, A. Way, SMT versus NMT: Preliminary comparisons for irish, in: Proceedings of the AMTA 2018 Workshop on Technologies for MT of Low Resource Languages (LoResMT 2018), Association for Machine Translation in the Americas, Boston, MA, 2018, pp. 12–20. URL: <https://www.aclweb.org/anthology/W18-2202>.
 - [619] I. Jauregi Unanue, L. Garmendia Arratibel, E. Zare Borzeshi, M. Piccardi, English-Basque statistical and neural machine translation, in: Proceedings of the 11th Language Resources and Evaluation Conference, European Language Resource Association, Miyazaki, Japan, 2018. URL: <https://www.aclweb.org/anthology/L18-1141>.
 - [620] A. K. Ojha, K. D. Chowdhury, C.-H. Liu, K. Saxena, The RGNLP machine translation systems for WAT 2018, arXiv preprint arXiv:1812.00798 (2018).
 - [621] Y. Jia, M. Carl, X. Wang, Post-editing neural machine translation versus phrase-based machine translation for English–Chinese, *Machine Translation* (2019) 1–21.
 - [622] P. Arthur, G. Neubig, S. Nakamura, Incorporating discrete translation lexicons into neural machine

- translation, in: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Austin, Texas, 2016, pp. 1557–1567. URL: <https://www.aclweb.org/anthology/D16-1162>. doi:10.18653/v1/D16-1162.
- [623] J. Zhang, C. Zong, Bridging neural machine translation and bilingual dictionaries, arXiv preprint arXiv:1610.07272 (2016).
 - [624] Y. Tang, F. Meng, Z. Lu, H. Li, P. L. Yu, Neural machine translation with external phrase memory, arXiv preprint arXiv:1606.01792 (2016).
 - [625] S. Vogel, H. Ney, C. Tillmann, HMM-based word alignment in statistical translation, in: COLING 1996 Volume 2: The 16th International Conference on Computational Linguistics, 1996. URL: <https://www.aclweb.org/anthology/C96-2141>.
 - [626] G. Neubig, M. Morishita, S. Nakamura, Neural reranking improves subjective quality of machine translation: NAIST at WAT2015, in: Proceedings of the 2nd Workshop on Asian Translation (WAT2015), Workshop on Asian Translation, Kyoto, Japan, 2015, pp. 35–41. URL: <https://www.aclweb.org/anthology/W15-5003>.
 - [627] R. Grundkiewicz, M. Junczys-Dowmunt, Near human-level performance in grammatical error correction with hybrid machine translation, in: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), Association for Computational Linguistics, New Orleans, Louisiana, 2018, pp. 284–290. URL: <https://www.aclweb.org/anthology/N18-2046>. doi:10.18653/v1/N18-2046.
 - [628] E. Avramidis, V. Macketanz, A. Burchardt, J. Helcl, H. Uszkoreit, Deeper machine translation and evaluation for German, in: Proceedings of the 2nd Deep Machine Translation Workshop, ÚFAL MFF UK, Lisbon, Portugal, 2016, pp. 29–38. URL: <https://www.aclweb.org/anthology/W16-6404>.
 - [629] B. Marie, A. Fujita, A smorgasbord of features to combine phrase-based and neural machine translation, in: Proceedings of the 13th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Papers), Association for Machine Translation in the Americas, Boston, MA, 2018, pp. 111–124. URL: <https://www.aclweb.org/anthology/W18-1811>.
 - [630] J. Zhang, M. Utiyama, E. Sumita, G. Neubig, S. Nakamura, Improving neural machine translation through phrase-based forced decoding, in: Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Asian Federation of Natural Language Processing, Taipei, Taiwan, 2017, pp. 152–162. URL: <https://www.aclweb.org/anthology/I17-1016>.
 - [631] F. Stahlberg, E. Hasler, B. Byrne, The edit distance transducer in action: The University of Cambridge English-German system at WMT16, in: Proceedings of the First Conference on Machine Translation, Association for Computational Linguistics, Berlin, Germany, 2016, pp. 377–384. URL: <https://www.aclweb.org/anthology/W16-2324>. doi:10.18653/v1/W16-2324.
 - [632] J. Niehues, E. Cho, T.-L. Ha, A. Waibel, Pre-translation for neural machine translation, in: Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers, The COLING 2016 Organizing Committee, Osaka, Japan, 2016, pp. 1828–1836. URL: <https://www.aclweb.org/anthology/C16-1172>.
 - [633] L. Zhou, W. Hu, J. Zhang, C. Zong, Neural system combination for machine translation, in: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 378–384. URL: <https://www.aclweb.org/anthology/P17-2060>. doi:10.18653/v1/P17-2060.
 - [634] J. Du, A. Way, Neural pre-translation for hybrid machine translation, in: Proceedings of MT Summit, volume 16, 2017, pp. 27–40.
 - [635] X. Wang, Z. Lu, Z. Tu, H. Li, D. Xiong, M. Zhang, Neural machine translation advised by statistical machine translation, in: Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, AAAI’17, AAAI Press, 2017, pp. 3330–3336. URL: <http://dl.acm.org/citation.cfm?id=3298023.3298052>.
 - [636] X. Wang, Z. Tu, M. Zhang, Incorporating statistical machine translation word knowledge into neural machine translation, IEEE/ACM Transactions on Audio, Speech, and Language Processing 26 (2018) 2255–2266.
 - [637] Z. Long, T. Utsuro, T. Mitsuhashi, M. Yamamoto, Translation of patent sentences with a large vocabulary of technical terms using neural machine translation, in: Proceedings of the 3rd Workshop on Asian Translation (WAT2016), The COLING 2016 Organizing Committee, Osaka, Japan, 2016, pp. 47–57. URL: <https://www.aclweb.org/anthology/W16-4602>.
 - [638] L. Dahlmann, E. Matusov, P. Petrushkov, S. Khadivi, Neural machine translation leveraging

phrase-based models in a hybrid search, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 1411–1420. URL: <https://www.aclweb.org/anthology/D17-1148>. doi:10.18653/v1/D17-1148.