



# Neural machine translation: A review of methods, resources, and tools

Zhixing Tan<sup>a,c,d</sup>, Shuo Wang<sup>a,c,d</sup>, Zonghan Yang<sup>a,c,d</sup>, Gang Chen<sup>a,c,d</sup>, Xuancheng Huang<sup>a,c,d</sup>,  
Maosong Sun<sup>a,c,d,e</sup>, Yang Liu<sup>a,b,c,d,e,\*</sup>

<sup>a</sup> Department of Computer Science and Technology, Tsinghua University, China

<sup>b</sup> Institute for AI Industry Research, Tsinghua University, China

<sup>c</sup> Institute for Artificial Intelligence, China

<sup>d</sup> Beijing National Research Center for Information Science and Technology, China

<sup>e</sup> Beijing Academy of Artificial Intelligence, China

## ARTICLE INFO

### Keywords:

neural Machine translation  
Attention mechanism  
Deep learning  
Natural language processing

## ABSTRACT

Machine translation (MT) is an important sub-field of natural language processing that aims to translate natural languages using computers. In recent years, end-to-end neural machine translation (NMT) has achieved great success and has become the new mainstream method in practical MT systems. In this article, we first provide a broad review of the methods for NMT and focus on methods relating to architectures, decoding, and data augmentation. Then we summarize the resources and tools that are useful for researchers. Finally, we conclude with a discussion of possible future research directions.

## 1. Introduction

Machine Translation (MT) is an important task that aims to translate natural language sentences using computers. The early approach to machine translation relies heavily on hand-crafted translation rules and linguistic knowledge. As natural languages are inherently complex, it is difficult to cover all language irregularities with manual translation rules. With the availability of large-scale parallel corpora, data-driven approaches that learn linguistic information from data have gained increasing attention. Unlike rule-based machine translation, Statistical Machine Translation (SMT) (Brown et al., 1990; Koehn et al., 2003) learns latent structures such as word alignments or phrases directly from parallel corpora. Incapable of modeling long-distance dependencies between words, the translation quality of SMT is far from satisfactory. With the breakthrough of deep learning, Neural Machine Translation (NMT) (Kalchbrenner and Blunsom, 2013; Cho et al., 2014a; Sutskever et al., 2014; Bahdanau et al., 2015) has emerged as a new paradigm and quickly replaced SMT as the mainstream approach to MT.

Neural machine translation is a radical departure from previous machine translation approaches. On the one hand, NMT employs continuous representations instead of discrete symbolic representations in SMT. On the other hand, NMT uses a single large neural network to model the entire translation process, freeing the need for excessive feature engineering. The training of NMT is end-to-end as opposed to separately

tuned components in SMT. Besides its simplicity, NMT has achieved state-of-the-art performance on various language pairs (Junczys-Dowmunt et al., 2016). In practice, NMT also becomes the key technology behind many commercial MT systems (Wu et al., 2016; Hassan et al., 2018).

As neural machine translation attracts much research interest and grows into an area with many research directions, we believe it is necessary to conduct a comprehensive review of NMT. In this work, we will give an overview of the key ideas and innovations behind NMT. We also summarize the resources and tools that are useful and easily accessible. We hope that by tracing the origins and evolution of NMT, we can stand on the shoulder of past studies, and gain insights into the future of NMT.

The remainder of this article is organized as follows: Section 2 will review the methods of NMT. We first introduce the basics of NMT, and then we selectively describe the recent progress of NMT. We focus on methods related to architectures, decoding, and data augmentation. Section 3 will summarize the resources such as parallel or monolingual corpora that are publicly available to researchers. Section 4 will describe tools that are useful for training and evaluating NMT models. Finally, we conclude and discuss future directions in Section 5.

\* Corresponding author. Department of Computer Science and Technology, Tsinghua University, China  
E-mail address: [liuyang2011@tsinghua.edu.cn](mailto:liuyang2011@tsinghua.edu.cn) (Y. Liu).

<https://doi.org/10.1016/j.aiopen.2020.11.001>

Received 10 October 2020; Accepted 20 November 2020

Available online 4 March 2021

2666-6510/© 2020 The Author(s). Published by Elsevier B.V. on behalf of KeAi Communications Co., Ltd. This is an open access article under the CC BY-NC-ND

license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

## 2. Methods

As a data-driven approach to machine translation, NMT also embraces the probabilistic framework. Mathematically speaking, the goal of NMT is to estimate an unknown conditional distribution  $P(y|x)$  given the dataset  $\mathcal{D}$ , where  $x$  and  $y$  are random variables representing source input and target output, respectively. We strive to answer the three basic questions of NMT:

- **Modeling.** How to design neural networks to model the conditional distribution?
- **Inference.** Given a source input, how to generate a translation sentence from the NMT model?
- **Learning.** How to effectively learn the parameters of NMT from data?

In 2.1, we first describe the basic methods of NMT for addressing the above three questions. We then dive into the details of NMT architectures in 2.2. We introduce non-autoregressive NMTs and bidirectional inference in 2.3, and discuss alternative training objectives and using monolingual data in 2.4 and 2.5.

Despite the great success, NMT is far from perfect. There are several theoretical and practical challenges faced by NMT. We survey the research progress of some important directions. We describe methods for open vocabulary in 2.6, prior knowledge integration in 2.7, and interpretability and robustness in 2.8.

### 2.1. Overview of NMT

#### 2.1.1. Modeling

Translation can be modeled at different levels, such as document-, paragraph-, or sentence-level. In this article, we focus on sentence-level translation. Besides, we also assume the input and output sentences are sequences. Thus the NMT model can be viewed as a *sequence-to-sequence* model. Assuming we are given a source sentence  $x = \{x_1, \dots, x_S\}$  and a target sentence  $y = \{y_1, \dots, y_T\}$ . By using the chain rule, the conditional distribution can factorize from left-to-right (L2R) as:

$$P(y|x) = \prod_{t=1}^T P(y_t|y_0, \dots, y_{t-1}, x). \quad (1)$$

NMT models which conform the Eq. (1) is referred as *L2R autoregressive NMT* (Kalchbrenner and Blunsom, 2013; Cho et al., 2014a; Sutskever et al., 2014; Bahdanau et al., 2015), for the prediction at time-step  $t$  is taken as an input at time-step  $t + 1$ .

Almost all neural machine translation models employ the *encoder-decoder framework* (Cho et al., 2014a). The encoder-decoder framework

consists of four basic components: the embedding layers, the encoder and decoder networks, and the classification layer. Fig. 1 shows a typical autoregressive NMT model using the encoder-decoder framework, which we shall use as an example. “<bos>” and “<eos>” are special symbols that mark the beginning and ending of a sentence, respectively.

The embedding layer embodies the concept of *continuous representation*. It maps a discrete symbols  $x_t$  into a continuous vector  $x_t \in \mathbb{R}^d$ , where  $d$  denotes the dimension of the vector. The embeddings are then fed into later layers for more finer-grained feature extraction.

The encoder network maps the source embeddings into hidden continuous representations. To learn expressive representations, the encoder must be able to model the ordering and complex dependencies that existed in the source language. Recurrent neural networks (RNN) are suitable choice for modeling variable-length sequences. With RNNs, the computation involves in encoder can be described as:

$$h_t = \text{RNN}_{\text{ENC}}(x_t, h_{t-1}). \quad (2)$$

By iteratively applying the state transition function  $\text{RNN}_{\text{ENC}}$  over input sequence, we can use the final state  $h_S$  as the representation for the entire source sentence, and then feed it to the decoder.

The decoder can be viewed as a language model conditioned on  $h_S$ . The decoder network extracts necessary information from the encoder output, and also models the long-distance dependencies between target words. Given the start symbol  $y_0 = \text{< bos >}$  and the initial state  $s_0 = h_S$ , the RNN decoder compresses the decoding history  $\{y_0, \dots, y_{t-1}\}$  into a state vector  $s_t \in \mathbb{R}^d$ :

$$s_t = \text{RNN}_{\text{DEC}}(y_{t-1}, s_{t-1}). \quad (3)$$

The classification layer predicts the distribution of target tokens. The classification layer is typically a linear layer with *softmax* activation function. Assuming the vocabulary of target language is  $\mathcal{V}$ , and  $|\mathcal{V}|$  is the size of the vocabulary. Given an decoder output  $s_t \in \mathbb{R}^d$ , the classification layer first maps  $h$  to a vector  $z$  in the vocabulary space  $\mathbb{R}^{|\mathcal{V}|}$  with the linear map. Then the softmax function is used to ensure the output vector is a valid probability:

$$\text{softmax}(z) = \frac{\exp(z)}{\sum_{i=1}^{|\mathcal{V}|} \exp(z_{[i]})}, \quad (4)$$

where we use  $z_{[i]}$  to denote the  $i$ -th component in  $z$ .

#### 2.1.2. Inference

Given an NMT model and a source sentence  $x$ , how to generate a translation from the model is an important problem. Ideally, we would like to find the target sentence  $y$  which maximizes the model prediction

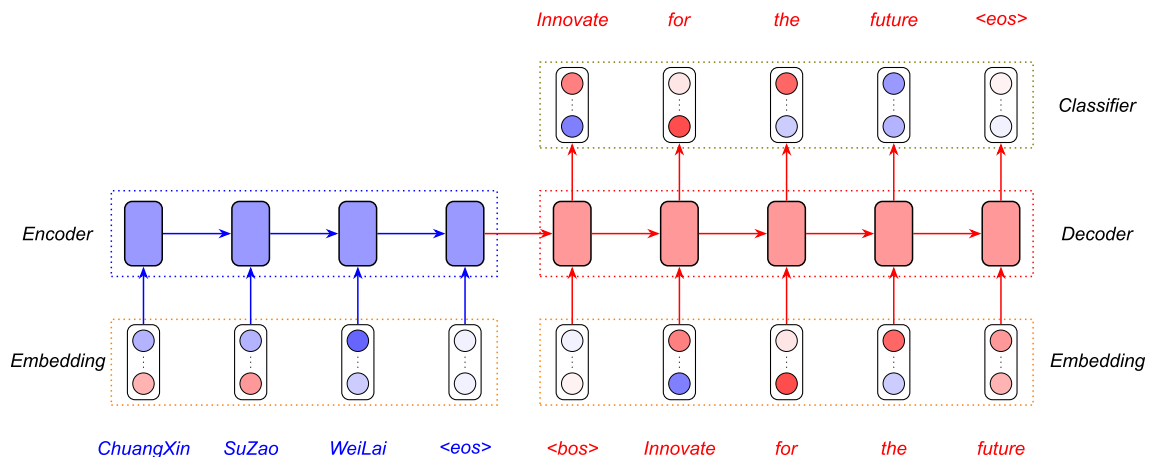


Fig. 1. An overview of the NMT architecture, which consists of embedding layers, a classification layer, an encoder network, and a decoder network. We use different colors to distinguish different languages.

$P(y|x=x;\theta)$  as the translation. However, due to the intractably large search space, it is impractical to find the translation with the highest probability. Therefore, NMT typically uses local search algorithms such as *greedy search* or *beam search* to find a local best translation.

Beam search is a classic local search algorithm which have been widely used in NMT. Previously, beam search have been successfully applied in SMT. The beam search algorithm keeps track of  $k$  states during the inference stage. Each state is a tuple  $\langle y_0 \dots y_t, v \rangle$ , where  $y_0 \dots y_t$  is a candidate translation, and  $v$  is the log-probability of the candidate. At each step, all the successors of all  $k$  states are generated, but only the top- $k$  successors selected. The algorithm usually terminates when the step exceed a pre-defined value or  $k$  full translation are found. It should be noted that the beam search will degrade into the greedy search if  $k = 1$ .

**Algorithm 1.** The beam search algorithm

---

**Algorithm 1:** The beam search algorithm

---

```

1  $t \leftarrow 1$  ;
2  $\mathcal{A} = \{ \langle \text{<bos>, } 0 \rangle \}$  ;  $\triangleright$  The set of alive candidates
3  $\mathcal{F} = \{ \}$  ;  $\triangleright$  The set of finished candidates
4 while  $t < \text{max\_length}$  do
5    $\mathcal{C} = \{ \}$  ;
6   for  $\langle y_0 \dots y_{t-1}, v \rangle \in \mathcal{A}$  do
7      $p \leftarrow \text{NMT}(y_0 \dots y_{t-1}, x)$  ;
8     for  $w \in \mathcal{V}$  do
9        $y_t \leftarrow w$  ;
10       $l \leftarrow \log(p[w])$  ;
11       $\mathcal{C} \leftarrow \mathcal{C} \cup \{ \langle y_0 \dots y_t, v + l \rangle \}$  ;
12    end
13  end
14   $\mathcal{C} \leftarrow \text{TopK}(\mathcal{C}, k)$  ;
15  for  $\langle y_0 \dots y_t, v \rangle \in \mathcal{C}$  do
16    if  $y_t == \text{<eos>}$  then
17       $\mathcal{F} \leftarrow \mathcal{F} \cup \{ \langle y_0 \dots y_t, v \rangle \}$  ;
18    else
19       $\mathcal{A} \leftarrow \mathcal{A} \cup \{ \langle y_0 \dots y_t, v \rangle \}$  ;
20    end
21  end
22   $t \leftarrow t + 1$  ;
23 end
24  $\langle y_0 \dots y_t, v \rangle \leftarrow \text{Top}(\mathcal{F})$  ;
25 return  $y_1 \dots y_t$ 
```

---

The pseudo-codes of the beam search algorithm are given in 1. We also give a running example of the algorithm in Fig. 2.

### 2.1.3. Training of NMT models

NMT typically uses maximum log-likelihood (MLE) as the training objective function, which is a commonly used method of estimating the parameters of a probability distribution. Formally, given the training set  $\mathcal{D} = \{ \langle x^{(s)}, y^{(s)} \rangle \}_{s=1}^S$ , the goal of training is to find a set of model parameters that maximize the log-likelihood on the training set:

$$\hat{\theta}_{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} \{ \mathcal{L}(\theta) \}, \quad (5)$$

where the log-likelihood is defined as

$$\mathcal{L}(\theta) = \sum_{s=1}^S \log P(y^{(s)} | x^{(s)}; \theta). \quad (6)$$

By the virtue of *back-propagation* algorithm, we can efficiently compute the gradient of  $\mathcal{L}$  with respect to  $\theta$ . The training of NMT models usually adopts *stochastic gradient search* (SGD) algorithm. Instead of computing gradients on the full training set, SGD computes the loss function and gradients on a *minibatch* of the training set. The plain SGD optimizer updates the parameters of an NMT model with the following rule:

$$\theta \leftarrow \theta - \alpha \nabla \mathcal{L}(\theta), \quad (7)$$

where  $\alpha$  is the *learning rate*. With well-chosen learning rate, the parameters of NMT are guaranteed to converge into a local optima. In practice, instead of plain SGD optimizer, adaptive learning rate optimizers such as Adam (Kingma and BaAdam, 2014) are found to greatly reduce the training time.

## 2.2. Architectures

### 2.2.1. Evolution of NMT architectures

Since 2013, there are attempts to build a pure neural MT. Early NMT architectures such as RCTM (Kalchbrenner and Blunsom, 2013), RNNencdec (Cho et al., 2014a), and Seq2Seq (Sutskever et al., 2014) adopt a *fixed-length* approach, where the size of source representation is fixed regardless the length of source sentences. These works typically use

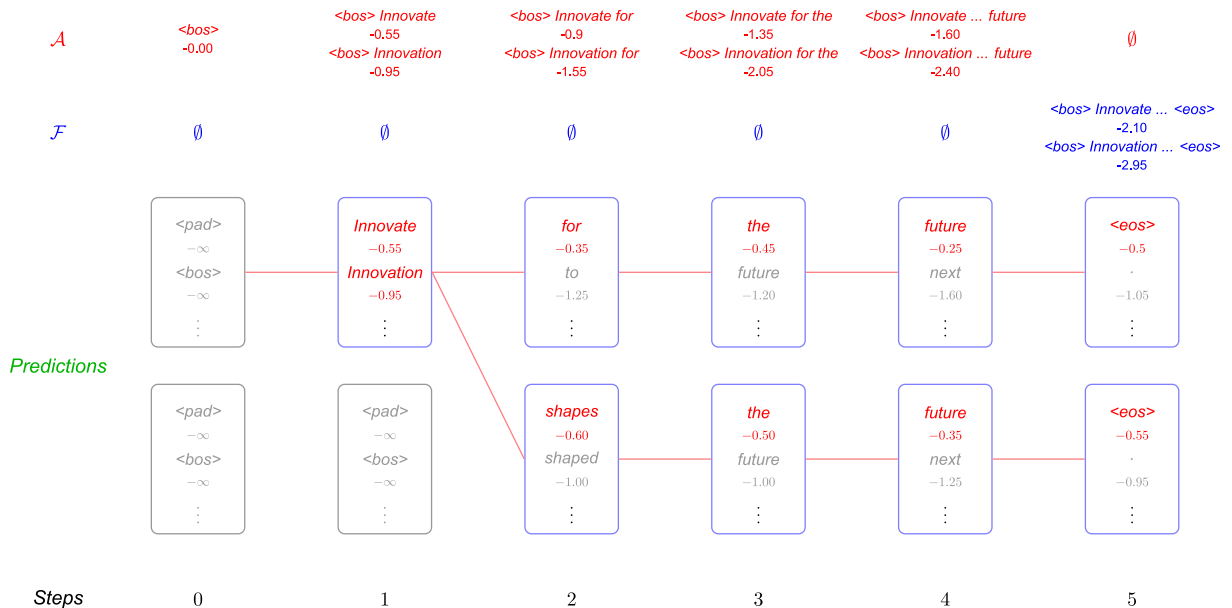


Fig. 2. A running example of the beam-search algorithm.

recurrent neural networks (RNN) as the decoder network for generating variable-length translation. However, it is found that the performance of this approach degrades as the length of the input sentence increases (Cho et al., 2014b). Two explanations can account for this phenomenon:

1. The fixed-length representations have become the bottleneck during the encoding process for long sentences (Cho et al., 2014a). As the encoder is forced to compress the entire source sentence into a set of fixed-length vectors, some important information may be lost in this process.
2. The longest path between the source words and target words is  $O(S + T)$ , and it is challenging for neural networks to learn long-term dependencies (Bengio et al., 1994). Sutskever et al. (2014) found that reverse the source sentence can significantly improve the performance of the fixed-length approach. By reversing the source sentence, the paths between the beginning words of source and target sentences are reduced, thus the optimization problem becomes easier.

Due to these limitations, later NMT architectures switch to *variable-length* source representations, where the length of source representations depends on the length of the source sentence. The RNNsearch architecture (Bahdanau et al., 2015) introduces *attention mechanism*, which is an important approach to implementing variable-length representations. Fig. 3 shows the comparison between fixed-length and variable-length approaches. By using the attention mechanism, the paths between any source and target words are within a constant length. As a result, the attention mechanism has eased optimization difficulty.

With the breakthrough of deep learning, NMT with deep neural networks have attracted much research interest. Seq2Seq (Sutskever et al., 2014) is the first architecture demonstrate the potential of deep NMT. Later architectures such as GNMT (Wu et al., 2016), ByteNet (Kalchbrenner et al., 2016), ConvSeq2Seq (Gehring et al., 2017), and Transformer (Vaswani et al., 2017) all use multi-layered neural networks. ByteNet and ConvSeq2Seq have replaced RNNs with convolutional neural networks (CNN) in their architectures while Transformer relies entirely on self-attention networks (SAN). Both CNNs and SANs can reduce the sequential operations involved in RNNs, and benefit from the parallel computation provided by modern devices such as GPU or TPU. Importantly, SAN can further reduce the longest path between two target tokens.

### 2.2.2. Attention mechanism

The introduction of attention mechanism (Bahdanau et al., 2015) is a milestone in NMT architecture research. The attention network computes the relevance of each value vector based on queries and keys. This can also be interpreted as a content-based addressing scheme (Graves et al., 2014). Formally, given a set of  $m$  query vectors  $\mathbf{Q} \in \mathbb{R}^{m \times d}$ , a set of  $n$  key vectors  $\mathbf{K} \in \mathbb{R}^{n \times d}$  and associated value vectors  $\mathbf{V} \in \mathbb{R}^{n \times d}$ , the computation of attention network involves two steps. The first step is to compute the relevance between keys and values, which is formally described as:

$$\mathbf{R} = \text{score}(\mathbf{Q}, \mathbf{K}), \quad (8)$$

where  $\text{score}$  is a scoring function which have several alternatives.  $\mathbf{R} \in \mathbb{R}^{m \times n}$  is a matrix storing the relevance score between each keys and values. The next step is compute the output vectors. For each query vector, the corresponding output vector is expressed as a weighted sum of value vectors:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}(\mathbf{R}) \cdot \mathbf{V}. \quad (9)$$

Fig. 4 depicts the two steps involved in the computation of attention mechanism.

Considering on the scoring function, the attention networks can be roughly classified into two categories: additive attention (Bahdanau et al., 2015) and dot-product attention (Luong et al., 2015). The additive attention models score through a feed-forward neural network:

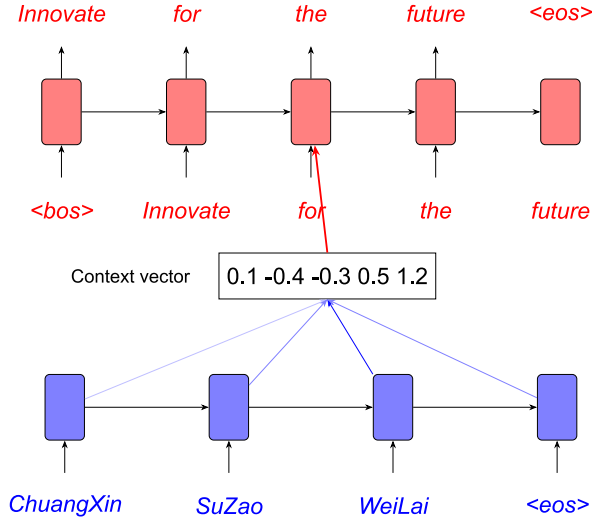


Fig. 3. At each decoding step, the attention mechanism dynamically generates a context vector based on the most relevant source representations for predicting the next target word.

$$\mathbf{R}_{[i,j]} = \mathbf{v}^\top \tanh(\mathbf{W}_s \mathbf{Q}_{[i]} + \mathbf{U}_s \mathbf{K}_{[j]}), \quad (10)$$

where  $\mathbf{W}_s \in \mathbb{R}^{d \times d}$ ,  $\mathbf{U}_s \in \mathbb{R}^{d \times d}$ , and  $\mathbf{v} \in \mathbb{R}^{d \times 1}$  are learnable parameters. On the other hand, the dot-product attention uses dot production to compute the matching score:

$$\mathbf{R}_{[i,j]} = \mathbf{Q}_{[i]}^\top \mathbf{K}_{[j]}. \quad (11)$$

In practice, the dot-product attention is much faster than the additive attention. However, the dot-product attention is found to be less stable than the additive attention when  $d$  is large (Vaswani et al., 2017). Vaswani et al. (2017) suspect that the dot-products grow large in magnitude for large values of  $d$ , which may resulting extremely small gradients caused by the softmax function. To remedy this issue, they propose to scale the dot-products by  $\frac{1}{\sqrt{d}}$ .

The attention mechanism is usually used as a part of the decoder network. Another type of attention network called self-attention network, is widely used in both the encoder and decoder of NMT. We shall describe self-attention and other variants of attention network later.

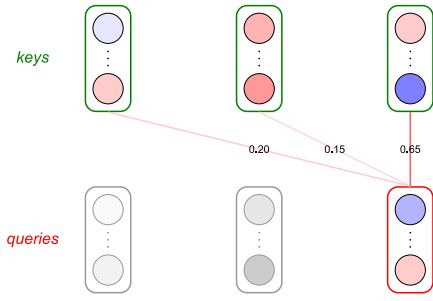
### 2.2.3. RNNs, CNNs, and SANs

The encoder and decoder are key components of NMT architectures. There are many methods to build powerful encoders and decoders, which can roughly divide into three categories: recurrent neural network (RNN) based methods, convolution neural network (CNN) based methods and self-attention network (SAN) based methods. There are several aspects we need to take into considerations for building an encoder and decoder:

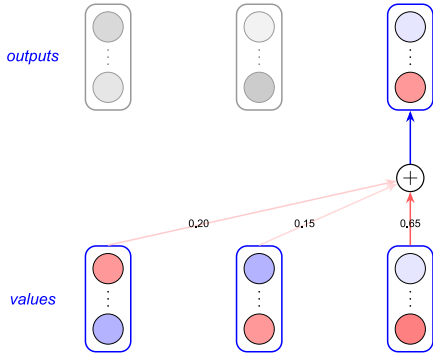
1. *Receptive field*. We hope each output produced by the encoder and decoder can potentially encode arbitrary information in the input sequence.
2. *Computational complexity*. It is desirable to a use network with lower computational complexity.
3. *Sequential operations*. Too many sequential operations preclude the parallel computation within the sequence.
4. *Position awareness*. The network should distinguish the ordering presents in the sequence.

Table 1 summarizes the computation as well as the above-mentioned aspects of typical RNN, CNN, and SAN (see Table 2).

Fig. 5 gives an overview of the ways of RNN, CNN, and SAN to encode sequences, respectively. As we can see in Fig. 5(a), RNNs are a family of sequential models that repeatedly apply the same state transition



(a) Given a query vector and key vectors, the attention network first computes a weight vector through the scoring function



(b) NMT with variable-length representation and attention mechanism.

Fig. 4. Detailed computations involved in the attention mechanism.

function to sequences. In theory, RNNs are among the most powerful family of neural networks (Siegelmann and Sontag, 1995). However, it suffers from severe vanishing and exploding gradient problem (Bengio et al., 1994) in practice. RNNs with gates, such as long short-term memory (LSTM) (Hochreiter and Schmidhuber, 1997) and gated recurrent unit (GRU) (Cho et al., 2014a) have been proposed to alleviate this problem. Another way to stabilize the training is to incorporate normalization layers, such as layer normalization (Ba et al., 2016).

In order to keep the auto-regressive property of NMT decoder during training, CNN and SAN further needs additional padding and masking to prevent the network from seeing future words. Fig. 6 shows padding and masking used in CNN and SAN.

Fig. 7 shows three extensions to RNNs that are widely used in NMT literature. Deep RNNs is one important way to increase the expressive power of RNNs. However, training deep neural networks is challenging because it also faces the vanishing and exploding gradient problem. There are many ways to construct deep RNNs, and the most popular one is by stacking multiple RNNs with *residual connections* (He et al., 2016a). The residual connection is an important method to construct deep neural networks. Residual connections use the identity mappings as the skip connections, which is formally described as:

$$\mathbf{y} = \mathbf{x} + f(\mathbf{x}), \quad (12)$$

where  $\mathbf{x}$ , and  $\mathbf{y}$  are input and output, respectively.  $f$  is the neural network. By using identity mappings, the gradient signal can directly propagate into lower layers. Bidirectional RNNs (Bahdanau et al., 2015) use two RNNs to process the same sequence in opposite directions, and concatenating the results of both RNNs to be the final output. In this way, each output of bidirectional RNNs encodes all the tokens in the sequence. An alternative to bidirectional RNNs is alternating RNNs (Zhou and Xu,

Table 1

Comparisons between different neural network layers. We use R.F. to denote the receptive field, S.O. to denote the number of sequential operations, and P.A. to denote the position awareness of the layer.  $t$  is the position in the sequence,  $l$  is the layer number. For CNN,  $k$  is the filter width and  $\mathbf{W}^{(i)}$  is the weight of the  $i$ -th filter.

Layer	Computation	R.F.	Complexity	S.O.	P.A.
RNN	$\mathbf{h}_{l,t} = \mathbf{W}\mathbf{h}_{l-1,t} + \mathbf{U}\mathbf{h}_{l,t-1}$	$\infty$	$O(n \cdot d^2)$	$O(n)$	Yes
CNN	$\mathbf{h}_{l,t} = \sum_{i=1}^k \mathbf{W}^{(i)} \mathbf{h}_{l-1,t+i} - \frac{k+1}{2}$	$k$	$O(k \cdot n \cdot d^2)$	$O(1)$	Yes
SAN	$\mathbf{h}_{l,t} = \sum_{i=1}^n \alpha_{l,t} \mathbf{h}_{l-1,i}$	$\infty$	$O(n^2 \cdot d)$	$O(1)$	No

Table 2

Comparison of fundamental architectures. V.R. denotes whether the architecture employs variable representation.  $\text{Path}_E$  denotes the longest path between the source and target tokens.  $\text{Path}_D$  denotes the longest path between two target tokens.

Model	Encoder	Decoder	Complexity	V.R.	$\text{Path}_E$	$\text{Path}_D$
RCTM 1 (Kalchbrenner and Blunsom, 2013)	CNN	RNN	$O(S^2 + T)$	No	$S$	$T$
RCTM 2 (Kalchbrenner and Blunsom, 2013)	CNN	RNN	$O(S^2 + T)$	Yes	$S$	$T$
RNN <sub>ENCDEC</sub> / Seq2Seq (Cho et al., 2014a; Sutskever et al., 2014)	RNN	RNN	$O(S + T)$	No	$S + T$	$T$
RNN <sub>SEARCH</sub> (Bahdanau et al., 2015)	RNN	RNN	$O(ST)$	Yes	1	$T$
ByteNet (Kalchbrenner et al., 2016)	CNN	CNN	$O(S + T)$	Yes	$c$	$c$
ConvSeq2Seq (Gehring et al., 2017)	CNN	CNN	$O(ST)$	Yes	1	$c$
TRANSFORMER (Vaswani et al., 2017)	SAN	SAN	$O(S^2 + ST + T^2)$	Yes	1	1

Table 3

Domain and language pairs provided by WMT20, IWSLT20, WAT20.

Workshop	Domain	Language Pair
WMT20	News	zh-en, cz-en, fr-de, de-en, iu-en, km-en, ja-en, ps-en, pl-en, ru-en, ta-en
	Biomedical	en-eu, en-zh, en-fr, en-de, en-it, en-pt, en-ru, en-es
IWSLT20	Chat	en-de
	TED Talks	en-de
	e-Commerce	zh-en, en-ru
	Open Domain	zh-ja
WAT20	Scientific Paper	en-ja, zh-ja
	Business Scene Dialogue	en-ja
	Patent	zh-ja, ko-ja, en-ja
	News	ja-en, ja-ru
	IT and Wikinews	hi-en, th-en, ms-en, id-en

2015), which consists of RNNs in opposite directions in adjacent layers.

Besides the difficulty training of RNNs, another major drawback of RNNs is that RNNs are sequential models in nature, which cannot benefit from the parallel computations provided by modern GPUs. CNNs and SANs, however, which fully exploit the parallel computation within sequences, are widely used in newer NMT architectures.



**Table 4**

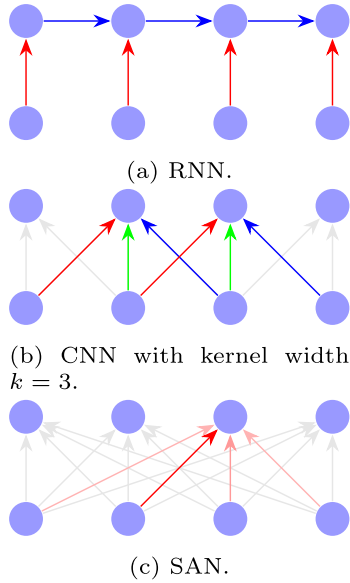
Number of sentences that available at OPUS for major languages to English.

Source	Fr-En	Es-En	De-En	Pt-En	Ru-En	Ar-En	Zh-En	Ja-En	Hi-En
OPUS (Tiedemann, 2016)	200.6M	172.0M	93.3M	77.7M	75.5M	69.2M	31.2M	6.2M	1.7M

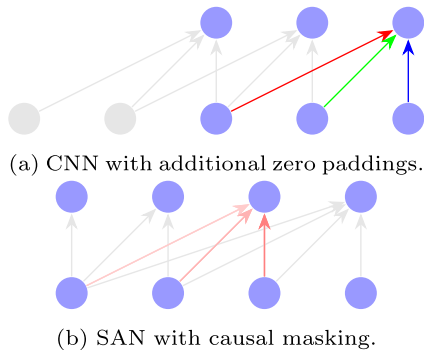
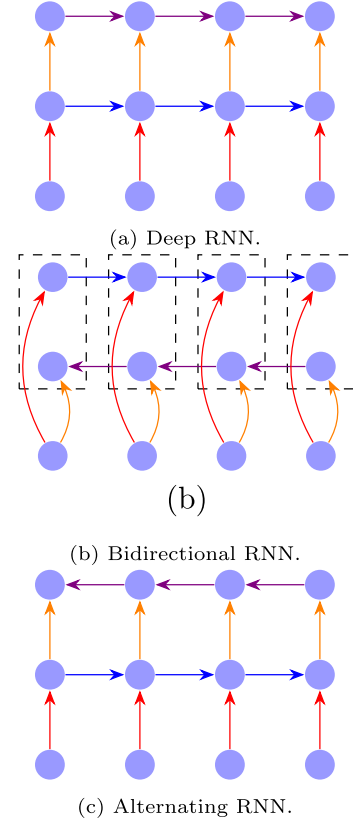
**Table 5**

Popular Open-source NMT toolkits on GitHub, the ordering is determined by the number of stars as the date of December 2020.

Name	Language	Framework	Status
TENSOR2TENSOR	Python	TensorFlow	Deprecated
FAIRSEQ	Python	PyTorch	Active
NMT	Python	TensorFlow	Deprecated
OPENNMT	Python/C++	PyTorch/TensorFlow	Active
SOCKEYE	Python	MXNet	Active
NEMATUS	Python	Tensorflow	Active
MARIAN	C++	–	Active
THUMT	Python	PyTorch/TensorFlow	Active
NMT-KERAS	Python	Keras	Active
NEURAL MONKEY	Python	TensorFlow	Active

**Fig. 5.** Overview of the computation diagram of RNN, CNN, and SAN. To be clarity, we use a node to denote the input or output vector of a specific layer.

Convolution neural network (CNN) was first introduced into NMT in 2013 (Kalchbrenner and Blunsom, 2013). However, it was not as successful as RNNs until 2017 (Gehring et al., 2017). The main obstacle for applying CNNs is its limited receptive field. Stacking  $L$  CNNs with kernel

**Fig. 6.** The computation of CNN and SAN during decoding.**Fig. 7.** Three extensions to RNNs.

width  $k$  can increase the receptive field from  $k$  to  $L \cdot (k - 1) + 1$ . The network needs to go deeper with large  $L$  and adopt large kernel size  $k$  to model long sentences. However, learning deep CNNs is challenging, and using large kernel size  $k$  may significantly increase the complexity and parameters involved in CNNs.

One solution to increase the receptive field without using a large  $k$  is through dilation (Kalchbrenner et al., 2016). Fig. 8 shows the comparison between plain CNN and dilated CNN. Plain CNN can be viewed as a special case of dilated CNN with a dilation rate  $r = 1$ . The computation of dilated CNN is mathematically formulated as:

$$\mathbf{h}_{l,t} = \sum_{i=1}^k \mathbf{W}^{(i)} \mathbf{h}_{l-1,t+\left(i-\left\lceil \frac{k+1}{2} \right\rceil\right) \times r} \quad (13)$$

Stacking  $L$  dilated CNNs whereby the dilation rates are doubled every layer, the receptive field increases to  $(2^L - 1) \cdot (k - 1) + 1$ . As a result, the receptive field grows exponentially with  $L$ , as opposed to linearly with  $L$  in plain CNN.

Another solution is to reduce the computations involved in CNN. Depthwise convolution (Kaiser et al., 2017) reduces the complexity from  $O(kd^2)$  to  $O(kd)$  by performing convolution independently over channels. Fig. 9 depicts the comparison between CNN and depthwise CNN. The output of the depthwise convolution layer is defined as:

$$\mathbf{h}_{l,t} = \sum_{i=1}^k \mathbf{w}^{(i)} \odot \mathbf{h}_{l-1,t+i-\left\lceil \frac{k+1}{2} \right\rceil} \quad (14)$$

where  $w^{(i)}$  is the  $i$ -th column of weight matrix  $\mathbf{W} \in \mathbb{R}$ . Lightweight convolution (Wu et al., 2019a) further reduces the number of parameters of depthwise convolution through weight sharing.

Self-attention network (SAN) (Vaswani et al., 2017) is a special case of attention network where the queries, keys, and values come from the same sequence. Similar to CNN, SAN is trivial to parallelize. Furthermore, Each output in SAN also has infinite receptive fields, which is the same with RNN. In SAN, the queries, keys, and values are typically obtained through a linear map of the input representations. The scaled dot-product self-attention mechanism can be formally described as:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}. \quad (15)$$

Multi-head attention (Vaswani et al., 2017) is an extended attention network with multiple parallel heads. Each head attends information from different subspace across value vectors. As a result, multi-head attention can perform more flexible transformations than the single-head attention. We give an illustration of multi-head attention in Fig. 10.

The major disadvantage of SAN network is that it ignores the ordering of words in the sequence. To remedy this, SAN needs additional position encoding to differentiate orders. Vaswani et al. (2017) proposed a sinusoid style position encoding, which is formulated as:

$$\text{timing}(t, 2i) = \sin(t / 10000^{2i/d}), \quad (16)$$

$$\text{timing}(t, 2i + 1) = \cos(t / 10000^{2i/d}), \quad (17)$$

where  $t$  is the position and  $i$  is the dimension index. Another popular way of position encoding is to learn an additional position embedding. Finally, the position encoding is added to each word representation, so the same words with different positions can have different representations.

#### 2.2.4. Comparison of fundamental architectures

We take the state-of-the-art Transformer architecture (Vaswani et al., 2017) as an example to put all things together. Fig. 11 shows the architecture of Transformer. The Transformer model relies solely on attention networks, with additional sinusoid-style position encoding added to input embedding. The transformer network consists of a stack of 6 encoder layers and 6 decoder layers. Each encoder layer contains two sub-layers whereas each decoder layer contains three sub-layers. To stabilize optimization, Transformer uses residual connection and layer

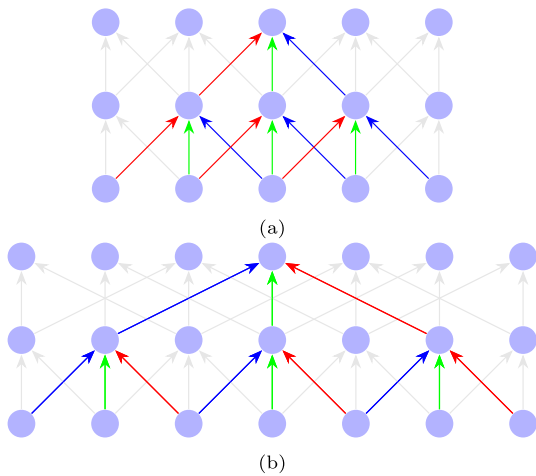


Fig. 8. Comparison between CNN and dilated CNN. (a) Two layers CNN with filter width  $k = 3$  for each layer. (b) Dilated CNN with filter width  $k = 3$  for all layers, dilation rate  $r = 1$  in layer 1 and  $r = 2$  in layer 2.

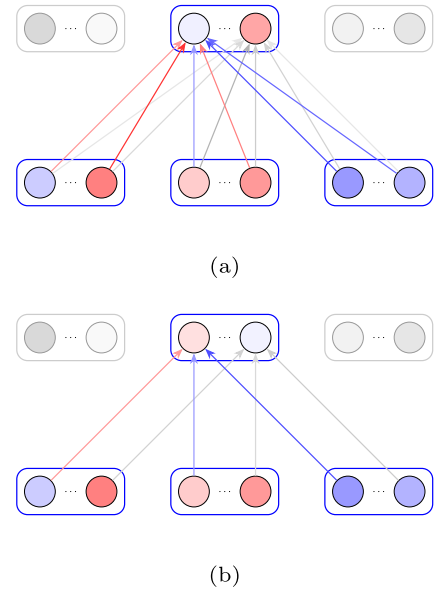


Fig. 9. Comparison between CNN and depthwise CNN. Each node in the graph represents a neuron instead of a vector. (a) Plain CNN. We highlight the computation of the first neuron in the output vector. (b) Depthwise CNN. Note that the connections are significantly reduced compared with plain CNN.

normalization in each sub-layer.

We summarize the comparison of fundamental NMT architectures in Table 2. We highlight several important aspects of these fundamental architectures.

#### 2.3. Bidirectional inference and non-autoregressive NMT

The dominate approach to NMT factorizes the conditional probability  $P(y|x)$  from left to right (L2R) auto-regressively. However, the factorization of the distribution is not unique. Researchers (Liu et al., 2016; Hoang et al., 2017; Zhang et al., 2018; Zhou et al., 2019a) have found that models with right-to-left (R2L) factorization are complementary to L2R models. The bidirectional inference is an approach to simultaneously generating translation with both L2R and R2L decoders. In addition to auto-regressive approaches where each output word on previously generated outputs, non-autoregressive NMTs (Gu et al., 2018) avoids this auto-regressive property and produces outputs in parallel, allowing much lower latency during inference.

##### 2.3.1. Bidirectional inference

Ignoring the future context is another obvious weakness of AR decoding. Thus, a natural idea is that the quality of translation will be improved if autoregressive models can “know” the future information. From this perspective, many approaches have been proposed to improve translation performance by exploring the future context. Some researchers proposed to model both past and future context (Zheng et al., 2018, 2019; Zhang et al., 2019a) and some others also found that L2R and R2L autoregressive models can generate complementary translations

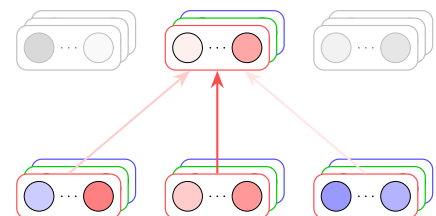


Fig. 10. An illustration of multi-head attention.

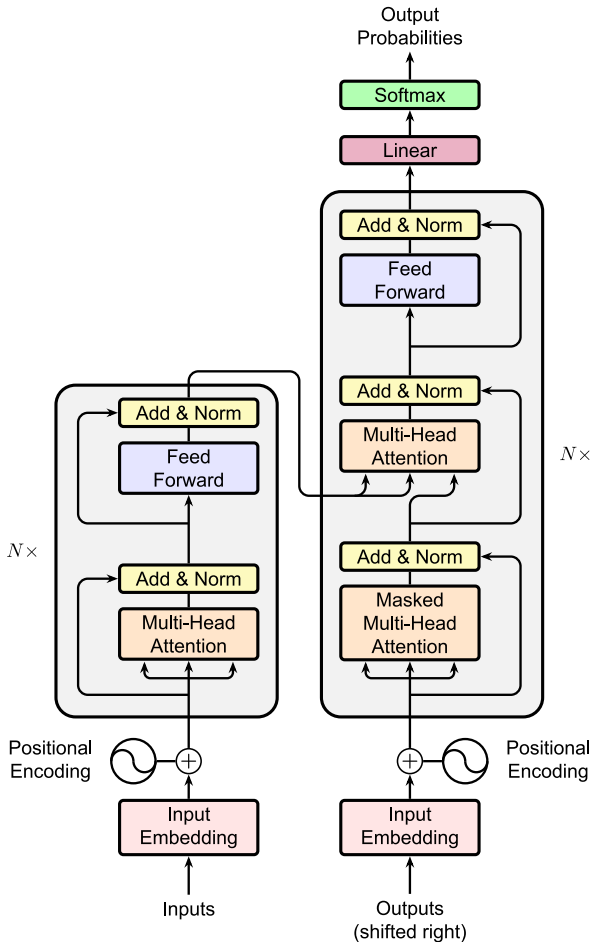


Fig. 11. The Transformer architecture.

(Liu et al., 2016; Hoang et al., 2017; Zhang et al., 2018; Zhou et al., 2019a). For instance, Zhou et al. (2019a) analyzed the translation accuracy of the first and last 4 tokens for L2R and R2L models, respectively. The statistical results show that, in Chinese-English translation, L2R performs better in the first 4 tokens while R2L translates better in the last 4 tokens.

Based on the findings mentioned above, a number of methods have been proposed to combine the advantages of L2R and R2L decoding. These approaches are collectively referred to as bidirectional decoding. Bidirectional decoding based methods can be mainly fall into four categories (Zhang and Zong, 2020): (1) agreement between L2R and R2L (Liu et al., 2016; Yang et al., 2018a; Zhang et al., 2019b), (2) rescore with bidirectional decoding (Liu et al., 2016; Sennrich et al., 2016a), (3) asynchronous bidirectional decoding (Zhang et al., 2018; Su et al., 2019), and (4) synchronous bidirectional decoding (Zhou et al., 2019a, 2019b; Zhang et al., 2020a).

Mathematically, the L2R translation order is rather arbitrary, and other arrangements such as R2L factorization are equally correct:

$$P(\mathbf{y}|\mathbf{x}) = \underbrace{\prod_{t=1}^T P(y_t|\mathbf{y}_{<t}, \mathbf{x})}_{\text{L2R model}} = \underbrace{\prod_{t=1}^T P(y_t|\mathbf{y}_{>t}, \mathbf{x})}_{\text{R2L model}}. \quad (18)$$

Based on this theoretical assumption, [Liu et al. (2016), Yang et al. (2018a), and Zhang et al. (2019b)] proposed joint training schemes in which each direction is used as a regularizer for the other direction. Empirical results show that these methods can lead to significant improvements compared with standard L2R and R2L models.

Another common scheme to combine L2R and R2L translations is rescoring (also known as reranking). A strong L2R model firstly produces

an  $n$ -best list of translations, and then an R2L model rescors each translation in the  $n$ -best list (Liu et al., 2016; Sennrich et al., 2016a, 2017). As the scores from L2R and R2L directions are based on complementary models, the quality of translation can be improved by rescoring. Recently, Zhang et al. (2018) introduced a new strategy to exploit both L2R and R2L models. They named this method asynchronous bidirectional decoding (ASBD), which first produces outputs (hidden states) by an R2L model and then uses these outputs to optimize the L2R model. ASBD can be done in three steps: The first step is to train a R2L model with bilingual corpora. The second step is to obtain outputs for each given source sentence using the trained R2L model. Finally, the output of R2L model is used as the additional context with the training data to train the L2R model. Thanks to incorporating the future information from the R2L model, the performance of L2R model can be substantially improved.

Although ASBD improves the quality of translation, it also incurs other problems. The L2R and R2L models are trained separately so that they have no chance to interact with each other. Besides, the L2R model translates source sentences based on the outputs of an R2L model, this degrades the efficiency of inference. To address these problems, Zhou et al. (2019a) further proposed a synchronous bidirectional decoding (SBD) method which generates translations using both L2R and R2L inference synchronously and interactively. Specifically, SBD uses a new synchronous attention model to allow both L2R and R2L models “communicating” with each other. As shown in Fig. 12(a), the dotted arrows illustrate interactions between L2R and R2L decoding. Zhou et al. (2019a) also designed a variant of the standard beam search algorithm to hold L2R and R2L decoding concurrently. The idea behind this algorithm is to maintain that each half beam contains L2R and R2L predictions, respectively. Empirical results show that SBD can significantly improve performance with a slight cost to decoding speed.

Mehri and Sigal (2018) proposed a novel middle-out decoder architecture that begins from an initial middle-word and simultaneously expands the sequence in both L2R and R2L directions. Zhou et al. (2019b) also proposed a similar method that allows L2R and R2L inferences to start concurrently from the left and right sides, respectively. Both L2R and R2L inferences terminate at the middle position. Extensive experiments demonstrate that this method can improve not only the accuracy of translation but also decoding efficiency.

### 2.3.2. Non-autoregressive NMTs

To reduce the latency during inference, Gu et al. (2018) first proposed the non-autoregressive NMT (NAT) to generate the target words in parallel. Formally, given the source sentence  $\mathbf{x}$ , the probability of the target sentence  $\mathbf{y}$  is modeled as follows:

$$P_{\mathcal{NAT}}(\mathbf{y}|\mathbf{x}; \theta) = P_L(T|\mathbf{x}; \theta) \cdot \prod_{t=1}^T P(y_t|\mathbf{x}; \theta), \quad (19)$$

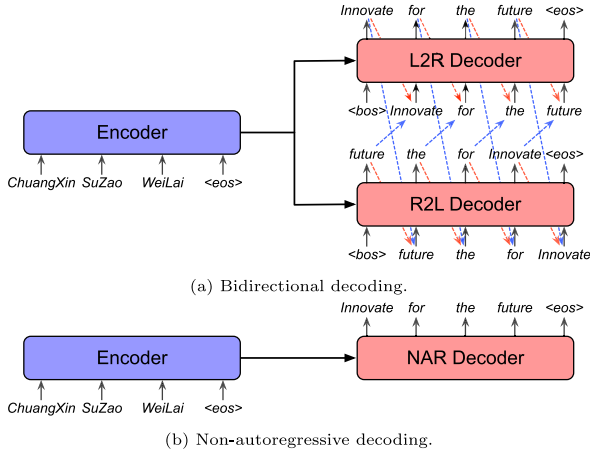
where  $P_{\mathcal{NAT}}(\mathbf{y}|\mathbf{x}; \theta)$  is the NAT model,  $P_L(T|\mathbf{x}; \theta)$  is a length sub-model to determine the length of target sentence, and  $\theta$  denotes the set of model parameters.

How to predict the length of target sentence (i.e.,  $P_L(T|\mathbf{x}; \theta)$  in Eq. (19)) is critical for NAT. Gu et al. (2018) proposed a fertility predictor to predict the length of translation. The fertility of a word in the source side determines how many target words it is aligned to. The fertility predictor can be denoted as

$$P_F(\mathbf{f}|\mathbf{x}; \theta) = \prod_{s=1}^S P(f_s|\mathbf{x}; \theta), \quad (20)$$

where  $\mathbf{f} = \{f_1, \dots, f_S\}$  is the fertility of the source sentence that consists of  $S$  words, and  $\theta$  is the set of parameters. At the training phase, the gold fertility of each sentence pair in the training data can be obtained by a word alignment system. At the inference phase, the length of the target sentence can be determined by the fertility predictor:





**Fig. 12.** Comparisons of different decoding strategies. (a) Bidirectional decoding: generates a sentence in both left-to-right (L2R) and right-to-left (R2L) directions; (b) Non-autoregressive (NAR) decoding: generates a sentence at one time.

$$\hat{T} = \sum_{s=1} \hat{f}_s, \quad (21)$$

$$\hat{f}_s = \operatorname{argmax}_{f_s} P(f_s | \mathbf{x}; \hat{\theta}), \quad (22)$$

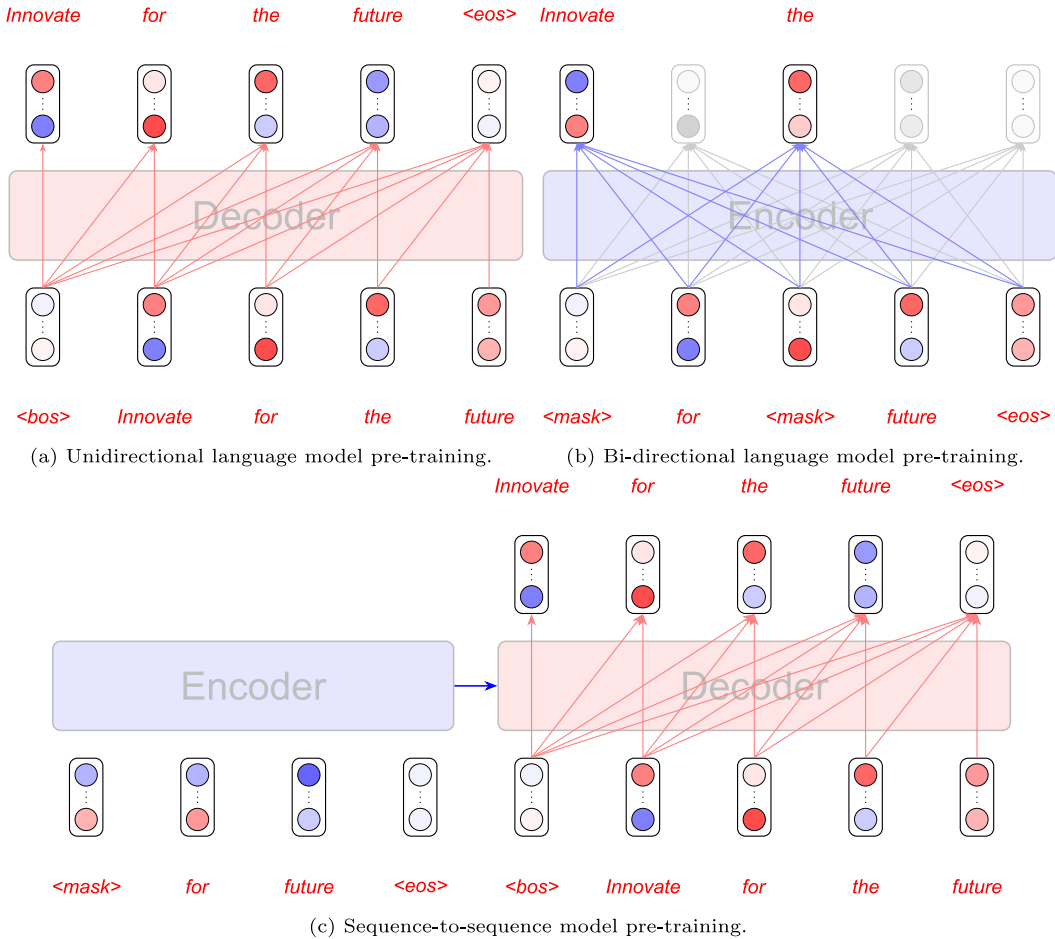
where  $\hat{T}$  is the number of words in the translation of the source sentence

$\mathbf{x}$ , and  $\hat{\theta}$  is the set of learned parameters.

Different from autoregressive NMT models that take the previous words (ie.,  $\mathbf{y}_{<t}$ ) as the input to predict the next target word  $y_t$ , NAT lacks such history information. Gu et al. (2018) also noticed that missing the input of the decoder can greatly impair translation quality. Thus, the authors proposed to copy each source token to the decoder, and the times each input token to be copied is its “fertility”. Gu et al. (2018) also used knowledge distillation (Kim and Rush, 2016), which employs strong autoregressive models as the “teachers” to improve the performance. Knowledge distillation has proven necessary for non-autoregressive translation (Zhou et al., 2019c; Gu et al., 2018; Lee et al., 2018; Libovický and Helcl, 2018; Ghazvininejad et al., 2019).

Despite the promising success of NAT, which can boost the decoding efficiency by about 15 times speedup compared with vanilla Transformer, NAT suffers from considerable quality degradation. Recently, many methods have been proposed to narrow the performance gap between non-autoregressive NMT and autoregressive NMT (Lee et al., 2018; Wang et al., 2018a, 2019a; Guo et al., 2019; Shao et al., 2019; Stern et al., 2019; Wei et al., 2019; Akoury et al., 2019; Ghazvininejad et al., 2019; Gu et al., 2019).

To take advantage of both autoregressive NMT and non-autoregressive NMT, Wang et al. (2018a) designed a semi-autoregressive Transformer (SAT) model. SAT keeps the autoregressive property in global but performs non-autoregressive translation in local. Specifically, SAT produces  $K$  sequential words per time-step independently to others. Consequently, SAT can balance autoregressive NMT ( $K = 1$ ) and non-autoregressive NMT ( $K = T$ ) by adjusting the value of  $K$ . Akoury et al. (2019) moved a further step to propose a syntactically supervised Transformer (SynST), which first autoregressively



**Fig. 13.** Three commonly used ways for pre-training.

predicts a chunked parse tree and then generates all words in one shot conditioned on the predicted parse.

A critical issue of NAT is that NAT copies the source words as the input of the decoder while ignores the difference between the source and target semantics. To address this problem, Guo et al. (2019) proposed to use a phrase table to covert source words to target words. They adopt a maximum match algorithm to greedily segment the source sentence into several phrases and then map these source phrases into target phrases by retrieving a pre-defined phrase table. Thanks to the enhanced decoder input, translation quality is significantly improved.

Inspired by the mask-predict task proposed by Devlin et al. (2019), Ghazvininejad et al. (2019) introduced a conditioned masked language model (CMLM) to generate translation by iterative refinement. CMLM trains the conditioned language model using a mask-predict manner and produces target sentences by iterative decoding during inference. Specifically, in the training phase, CMLM first randomly masks the words in the target sentence and then predicts these masked words. In the inference, CMLM generates the entire target sentence in a preset number of decoding iteration  $N$ . At iteration  $n \in [1, N]$ , the decoder input is the entire target sentence with  $T - \frac{T(N-t+1)}{N}$  words masked. The decoding process starts with a fully-masked target sentence and the words with the lowest prediction probabilities will be masked. With a proper number of decoding iteration, CMLM can effectively close the gap with fully autoregressive models and maintain the decoding efficiency.

#### 2.4. Alternative training objectives

NMT trained with maximum likelihood estimation or MLE have achieved state-of-the-art results on various language pairs (Junczys-Dowmunt et al., 2016). Despite the remarkable success, Ranzato et al. (2015) indicate two drawbacks of MLE for NMT. First, NMT models are not exposed to their errors during training, which is referred to as the *exposure bias* problem. Second, MLE is defined at word-level rather than sentence-level. Due to these limitations, researchers have investigated several alternative objectives.

Shen et al. (2016) proposed minimum risk training (MRT) to alleviate the problem. In MRT, the risk is defined as the expected loss with respect to the posterior distribution:

$$\mathcal{L}(\theta) = \sum_{s=1} \mathbb{E}_{\mathbf{y}|\mathbf{x}^{(s)}, \theta} [\Delta(\mathbf{y}, \mathbf{y}^{(s)})] \quad (23)$$

$$= \sum_{s=1} \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{x}^{(s)})} P(\mathbf{y}|\mathbf{x}^{(s)}; \theta) \Delta(\mathbf{y}, \mathbf{y}^{(s)}), \quad (24)$$

where  $\mathcal{Y}(\mathbf{x}^{(s)})$  is a set of all possible candidate translations for  $\mathbf{x}^{(s)}$ ;  $\Delta(\mathbf{y}, \mathbf{y}^{(s)})$  measures the difference between model prediction and gold-standard. Shen et al. (2016) indicate three advantages for MRT over MLE. Firstly, MRT direct optimize NMT with respect to evaluation metrics. Secondly, MRT can incorporate with arbitrary loss functions. Finally, MRT is transparent to architectures and can be applied to any end-to-end NMT systems. MRT achieves significant performance improvements than MLE training for RNNSearch. Reinforcement learning adopts a similar way as MRT does, which is comprehensively studied in Ranzato et al. (2016), Wu et al. (2018)]. However, recent literature (Choshen et al., 2020) has also pointed out the weakness of reinforcement learning for NMT, including discussion about authentic optimization goals and difficulty in convergence.

Efforts on improving training objectives reveal the art of translating motivation into functions and rewrite the conventional loss function with them or integrating them into it as regularizers. A collection of classical structured prediction losses are reviewed and compared in Edunov et al. (2018a), including MLE, sequence-level MLE, MRT, and max-margin learning. Yang et al. (2019a) leveraged the idea of max-margin learning in reducing word omission errors in NMT. They artificially constructed negative examples by omitting words in target reference

sentences, forcing the NMT model to assign a higher probability to a ground-truth translation and a lower probability to an erroneous translation. Wieting et al. (2019) aimed at improving the semantic similarity between ground-truth references and translation outputs from NMT systems. They proposed to use a margin-based loss as an alternative reward function, encouraging NMT models to output semantically correct hypotheses even if they mismatch with the reference in the lexicon. Chen et al. (2019) aimed at improving model capability of capturing long-range semantic structure. They proposed to explicitly model the source-target alignment with optimal transport (OT), and couple the OT loss with the MLE loss function. Kumar and Tsvetkov (2019) aimed at improving model efficiency and reducing the memory footprint of NMT models. Observing that the softmax layer usually takes considerable memory usage and the longest computation time, they proposed to replace the softmax layer with a continuous embedding layer, using Von Mises-Fisher distribution to implement soft ranking as softmax layer functions. As a result, the novel probabilistic loss enables NMT models to train much faster and handle very large vocabularies.

#### 2.5. Using monolingual data and unsupervised NMT

The amount of parallel data significantly affects the training of parameters as NMT is found to be data-hungry (Zoph et al., 2016). Unfortunately, large-scale parallel corpora are not available for the vast majority of language pairs. In contrast, monolingual corpora are abundant and much easy to obtain. As a result, it is important to augment the training set with monolingual data.

##### 2.5.1. Using monolingual data

As NMT is trained in an end-to-end way, it raises the difficulties in taking advantage of monolingual data. In the past few years researchers have proposed various methods to make use of the source- and target-side monolingual data in neural machine translation.

For target-side monolingual data, early attempts try to incorporate a language model trained on large-scale monolingual data into NMT. Gulcehre et al. (2017) proposed two ways to integrate a language model. One way is called *shallow fusion*, which uses a language model during decoding to rescore the  $n$ -base list. Another way is called *deep fusion*, which combines the decoder and language model with a controller mechanism. However, the improvements of these approaches are limited.

Another way to use target-side monolingual data is called *Back-translation* (BT) (Sennrich et al., 2016b). BT can make use of target-side monolingual data without changing the architecture of NMT. In Sennrich et al. (2016b), they first trained a target-to-source translation model using the parallel corpus. Then, the target-side monolingual data are used to build a synthetic parallel corpus, whose source sides are generated by the target-to-source translation model. Finally, the concatenation of parallel corpus and synthetic parallel corpus is used to learn a source-to-target translation model. Although the architecture and decoding algorithm is kept unchanged, the monolingual data is fully utilized to improve the translation quality. The authors attributed the effectiveness of using monolingual data to domain adaptation effects, reductions of overfitting, and improved fluency. BT has shown to be the most simple and effective method to leverage target-side monolingual data (Sennrich et al., 2016b; Poncelas et al., 2018). It is especially useful when only a small number of parallel data is available (Karakanta et al., 2018). Imamura et al. (2018) found that the diversities of source sentences affect the performance of BT. In the meantime, Edunov et al. (2018b) analyzed BT extensively and showed that noised-BT, which builds a synthetic corpus by sampled source sentences or noised output of beam-search, leads to higher accuracy. Caswell et al. (2019) investigated the role of noise in noised-BT. They revealed that the noises work in a way of making the model be able to distinguish the synthetic data and genuine data. The model can further take advantages of helpful signal and ignore harmful signal. As a result, they proposed a simple method called tagged-BT, which appends a preceding tag (e.g., <BT>) to every

synthetic source sentence. Wang et al. (2019b) proposed to consider uncertainty-based confidence to help NMT models distinguish synthetic data from authentic data.

Besides target-side monolingual data, source-side monolingual data are also important resources to improve the translation quality of semi-supervised neural machine translation. Zhang and Zong (2016) explored two ways to leverage source-side monolingual data. The former one is knowledge distillation (also called self-training), which utilizes the source-to-target translation model to build a synthetic parallel corpus. The latter is multi-task learning that simultaneously learns translation and source sentence reordering tasks.

There are many works to make use of both source- and target-side monolingual data. Hoang et al. (2018) found that the translation quality of the target-to-source model in BT matters and then proposed iterative back-translation, making the source-to-target and target-to-source to enhance each other iteratively. Cheng et al. (2016) presented an approach to train a bidirectional neural machine translation model, which introduced autoencoders on the monolingual corpora with source-to-target and target-to-source translation models as encoders and decoders by appending a reconstruction term to the training objective. He et al. (2016b) proposed a dual-learning mechanism, which utilized reinforcement learning to make the source-to-target and target-to-source model to teach each other with the help of source- and target-side language models. Zheng et al. (2020) proposed a mirror-generative NMT model to integrate source-to-target and target-to-source NMT models and both-side language models, which can learn from monolingual data naturally.

Pre-training is an alternative way to utilize monolingual data for NMT, which is shown to be beneficial by further combining with back-translation in the supervised and unsupervised NMT scenario (Song et al., 2019; Liu et al., 2020). Recently, pre-training has attracted tremendous attention because of its effectiveness on low-resource language understanding and language generation tasks (Peters et al., 2018; Radford et al., 2019; Devlin et al., 2019). Researchers found that models trained on large-scale monolingual data can learn linguistics knowledge (Clark et al., 2019). These knowledge can be transferred into downstream tasks by initializing the task-oriented models with the pre-trained weights. Language modeling is a commonly used pre-training method. The drawback of standard language modeling is that it is unidirectional, which may be sub-optimal as a pre-training technique. Devlin et al. (2019) proposed a masked pre-training language model (MLM) objective, which allows the model to make full use of context at the price of losing the ability to generate sequences. Combining language modeling and masking with sequence-to-sequence models, however, do not suffer from these limitations (Song et al., 2019; Lewis et al., 2019; Liu et al., 2020). Fig. 13 summarizes the three commonly used ways for pre-training.

Edunov et al. (2019) fed the output representations of ELMO (Peters et al., 2018) to the encoder of NMT. Zhu et al. (2020) proposed to fuse extracted representations into each layer of encoder and decoder through attention mechanism. Song et al. (2019) proposed to pre-train a sequence-to-sequence model first, and then finetune the pre-trained model on translation task directly. BART (Lewis et al., 2019) took various noising method to pre-train a denoising sequence-to-sequence model and then finetune the model with an additional encoder that replaces the word embeddings of the pre-trained encoder. Liu et al. (2019a) proposed mBART which is trained by applying BART to large-scale monolingual data across many languages.

### 2.5.2. Unsupervised NMT

Due to insufficient parallel corpus, it is not feasible to use supervised methods to train an NMT model on many language pairs. Unsupervised neural machine translation aims to obtain a translation model without using parallel data. Apparently, unsupervised machine translation is much more difficult than the supervised and semi-supervised settings.

Unsupervised neural machine translation is composed of three parts.

First, by the virtue of recent advances on unsupervised cross-lingual embeddings (Zhang et al., 2017a; Artetxe et al., 2017a) and word-by-word translation systems (Conneau et al., 2017), the unsupervised translation models can be initialized by weak translation models with fundamental cross-lingual information. Second, denoising autoencoders (Vincent et al., 2008) are used to embed the sentences into dense latent representations. The sentences of different languages are assumed to be embedded into the same latent space so that the latent representations of source sentences can be decoded into the target language. Third, iterative back-translation is used to strengthen the source-to-target and target-to-source translation models. Lample et al. (2017) and Artetxe et al. (2017b) first successfully built an unsupervised NMT system as described above. Specifically, Lample et al. (2017) utilizes a discriminator to force the encoder to embed sentences of each language to the same latent space.

While Lample et al. (2017) used a shared encoder and a shared decoder, Artetxe et al. (2017b) adopt a shared encoder but two separate decoder approach. Yang et al. (2018b) conjectured that sharing of the encoder and decoder between two languages may lose their language characteristics. Therefore they proposed leveraging two separate encoders with some shared layers and using two different GANs to restrict the latent representations. Artetxe et al. (2018) and Lample et al. (2018) found that an unsupervised statistical machine translation system with iterative back-translation can easily outperform the unsupervised NMT counterpart. Lample et al. (2018) summarized that initialization, language modeling, and iterative back-translation are three principles in fully unsupervised MT and they further found that combining unsupervised SMT and unsupervised NMT can reach better performances.

Ren et al. (2019) suggested that the noises and errors existed in pseudo-data can be accumulated and hinder the improvements during iterative back-translations. Therefore, they proposed to use SMT which is less sensitive to noises as posterior regularizations to unsupervised NMT. As the unsupervised NMT is usually initialized by unsupervised bilingual word embeddings (UBWE), Sun et al. (2019) proposed to utilize UBWE agreement to enhance unsupervised NMT. Wu et al. (2019b) considered that pseudo sentences predicted by weak unsupervised MT systems are usually of low quality. To alleviate this issue, they proposed an extract-edit approach, which is an alternative to back-translation. First, they extracted the most relevant target sentences from target monolingual data given the source sentence. Then, extracted target sentences were edited to be aligned with the source sentences. This method makes it possible to use real sentence pairs to train the unsupervised NMT system. Ren et al. (2020) also proposed a similar retrieve-and-rewrite method to initialize an unsupervised SMT system. Artetxe et al. (2019) improved unsupervised SMT by exploiting subword information, developing a theoretically well-founded unsupervised tuning method, and incorporating a joint refinement procedure. Finally, they utilized the improved unsupervised SMT to initialize NMT model and get state-of-the-art results. As a unique method to utilize monolingual data, cross-lingual pre-trained models are used by Lample and Conneau (2019) to initialize unsupervised MT systems.

### 2.6. Open vocabulary

NMT typically operates with a fixed vocabulary. Due to practical reasons such as computational concerns and memory constraints, the vocabulary size of NMT models often ranges from 30k to 50k. For word-level NMT, the limited size of vocabulary results in a large number of unknown words. Therefore, word-level NMT is unable to translate these words and performs poorly in open-vocabulary settings (Sutskever et al., 2014; Bahdanau et al., 2015).

Although word-level NMT is unable to translate out-of-vocabulary words, character-level NMT do not have this problem. By splitting words into characters, the vocabulary size is much smaller and every rare word can be represented. Chung et al. (2016) found that the NMT model with subword-level encoder and character-level decoder can also work

well. Lee et al. (2017) introduced a fully character-level NMT with convolutional network and found that character-to-character NMT is suitable in many-to-one multilingual setting. Luong and Manning (2016) built hybrid systems that translate mostly at the word level and consult the character components for rare words. Passban et al. (2018) proposed an extension to the model of Chung et al. (2016), which works at the character level and boosts the decoder with target-side morphological information. Chen et al. (2018) proposed an NMT model at different levels of granularity with a multi-level attention. Gao et al. (2020) found that self-attention performs very well on character-level translation.

Character-level NMT also has its imperfection, splitting words into characters results in longer sequences in which each symbol contains less information, creating both modeling and computational challenges (Cherry et al., 2018). Other than word-level and character-level methods, subword-level method is another choice to model input and output sentences. Sennrich et al. (2016c) first adapted byte-pair-encoding (BPE) to word segmentation task, which is a simple but effective method. BPE making the NMT model capable of open-vocabulary translation by encoding rare and unknown words as sequences of subword units. This method reaches a compromise between vocabulary size and sequence length with stabilized better performance over word- and character-level methods. Moreover, it is an unsupervised method with few hyper-parameters, making it the most commonly used method for word segmentation for neural machine translation and text generation. Kudo (2018) presented a simple regularization method, namely subword regularization, to improve the robustness of subword-level NMT. Provilkov et al. (2019) introduced BPE-dropout to regularize the subword segmentation algorithm BPE, which is more compatible with conventional BPE than the method proposed by Kudo (2018). Wang et al. (2020) investigated byte-level subwords, specifically byte-level BPE (BBPE), which is more efficient than using pure bytes only.

## 2.7. Prior knowledge integration

Prior knowledge such as parse tree has been shown to be effective to improve SMT (Liu et al., 2006). Several works find that integrating prior knowledge can also improve the translation performance of NMT models.

One line of studies focus on inducing lexical knowledge into NMT models. Zhang et al. (2017b) proposed a general framework that can integrate prior knowledge into NMT models through posterior regularization and found that bilingual dictionary is useful to improve NMT models. Morishita et al. (2018) found that feeding hierarchical subword units to different modules of NMT models can also improve the translation quality. Liu et al. (2019b) proposed a novel shared-private word embedding to capture the relationship of different words for NMT models. Chen et al. (2020) distinguished content words and functional words depending on the term frequency inverse document frequency (i.e., *TF-IDF*) and then added an additional encoder and an additional loss for content words. Weller-Di Marco and Fraser (Weller-Di Marco and Fraser, 2020) studied strategies to model word formation in NMT to explicitly model fusional morphology.

Modeling the source-side syntactic structure has also drawn a lot of attention. Eriguchi et al. (2016) extended NMT models to an end-to-end syntactic model, where the decoder is softly aligned with phrases at the source side when generating a target word. Sennrich and Haddow (2016) explored external linguistic information such as lemmas, morphological features, POS tags and dependency labels to improve translation quality. Hao et al. (2019) presented a multi-granularity self-attention mechanism to model phrases which are extracted by syntactic trees. Bugliarello and Okazaki (2020) proposed the Parent-Scaled Self-Attention to incorporate dependency tree to capture the syntactic knowledge of the source sentence. There are also some works that use multi-task training to learn source-side syntactic knowledge, in which the encoder of a NMT model is trained to perform POS tagging or syntactic parsing (Eriguchi et al., 2017; Baniata et al., 2018).

Another line of studies directly model the target-side syntactic

structures (Gü et al., 2018; Wang et al., 2018b; Wu et al., 2017; Aharoni and Goldberg, 2017; Bastings et al., 2017; Li et al., 2018; Yang et al., 2019b, 2020). Aharoni and Goldberg (2017) trained a end-to-end model to directly translate source sentences into constituency trees. Similar approaches are proposed to use two neural models to generate the target sentence and its corresponding tree structure (Wang et al., 2018b; Wu et al., 2017). Gü et al. (2018) proposed to use a single model to perform translation and parsing at the same time. Yang et al. (2019b) introduced a latent variable model to capture the co-dependence between syntax and semantics. Yang et al. (2020) trained a neural model to predict the soft template of the target sentence conditioning only on the source sentence and then incorporated the predicted template into the NMT model via a separate template encoder.

## 2.8. Interpretability and robustness

Despite the remarkable progress, it is hard to interpret the internal workings of NMT models. All internal information in NMT is represented as high-dimensional real-valued vectors or matrices. Therefore, it is challenging to associate these hidden states with language structures. The lack of interpretability has made it very difficult for researchers to understand the translation process of NMT models.

In addition to interpretability, the lack of robustness is a severe challenge for NMT systems as well. With small perturbations in source inputs (also referred to as *adversarial examples*), the translations of NMT models may lead to significant erroneous changes (Belinkov and Bisk, 2018; Cheng et al., 2019). The lack of robustness of NMT limits its application on tasks that require robust performance on noisy inputs. Therefore, improving the robustness of NMT has gained increasing attention in the NMT community.

### 2.8.1. Interpretability

Efforts have been devoted to improving the interpretability of NMT systems in recent works. Ding et al. (2017) proposed to visualize the internal workings of the RNNSearch (Bahdanau et al., 2015) architecture. With layer-wise relevance propagation (Bach et al., 2015), they computed and visualized the contribution of each contextual word to arbitrary hidden states in RNNSearch. Bau et al. (2019) share similar motivations with Ding et al. (2017). Their basic assumption is that the same neuron in different NMT models captures similar syntactic and semantic information. They proposed to use several types of correlation coefficients to measure the importance of each neuron. As a result, by identifying important neurons and controlling their activation, the translation process of NMT systems can be controlled. Strobel et al. (2019) also put effort into visualizing the working process of RNNSearch. The highlights of their work lie in the utilization of training data. When an NMT system decodes some words, their visualization system provides the most relevant training corpora by using the nearest neighbor search. In case of translation errors, the system can locate the erroneous outputs directly in the training set by showing its origin cause. As a result, this function provides better assistants and makes it easy for developers to adjust the model and the data.

With the tremendous success of the Transformer architecture (Vaswani et al., 2017), the NMT community have shown increasing interest in understanding and interpreting Transformer. He et al. (2019) generalized the idea of layer-wise relevance to word importance by attributing the NMT output to every input word through a gradient-based method. The calculated word importance illustrates the influence of each source words, which also serves as an implication of under-translation errors. Raganato and Tiedemann (2018) analyzed the internal representations of Transformer encoder. Utilizing the attention weights in each layer, they extract relation among each word in the source sentence. They designed four types of probing tasks to analyze the syntactic and semantic information encoded by each layer representation and test their transferability. Voita et al. (2019) also proposed to analyze the bottom-up evolution of representations in Transformer with canonical correlation



analysis (CCA). By estimating mutual information, they studied how information flows in Transformer. Stahlberg et al. (2018) proposed an operation sequence model to interpret NMT. Based on the translation outputted by the Transformer system, they proposed explicit modeling of the word reordering process and provided explicit word alignment between the reordered target-side sentence and the source sentence. As a result, one can track the reordering process of each word's information as they are explicitly aligned with the source side. Recent work (Yun et al., 2020) also provided a theoretical understanding of Transformer by proving that Transformer networks are universal approximators of sequence-to-sequence functions.

### 2.8.2. Robustness

Belinkov and Bisk (2018) first investigated the robustness of NMT. They pointed out that both synthetic and natural noise can severely harm the performance of NMT models. They experimented with four types of synthetic noise and leveraged structure-invariant representation and adversarial training to improve the robustness of NMT. Similarly, Zhao et al. (2018) proposed to map the input sentence to a latent space with generative adversarial networks (GAN) and search for adversarial examples in that space. Their approach can produce semantically and syntactically coherent sentences that have negative impacts on the performance of NMT models.

Ribeiro et al. (2018) proposed semantic-preserving adversarial rules to explicitly induce adversarial examples. This approach provides a better guarantee for the adversarial examples to satisfy semantically equivalence property. Cheng et al. (2018) proposed two types of approaches to generating adversarial examples by perturbing the source sentence or the internal representation of the encoder. By integrating the effect of adversarial examples into the loss function, the robustness of neural machine translation is improved by adversarial training.

Ebrahimi et al. (2018) proposed a character-level white-box attack for character-level NMT. They proposed to model the operations of character insertion, deletion, and swapping with vector computations so that the generation of adversarial examples can be formulated with differentiable string-edit operations. Liu et al. (2019a) proposed to jointly utilize textual and phonetic embedding in NMT to improve robustness. They found that to train a more robust model, more weights should be put on the phonetic rather than textual information. Cheng et al. (2019) proposed doubly adversarial inputs to improve the robustness of NMT. Concretely, they proposed to both attack the translation model with adversarial source examples and defend the translation model with adversarial target inputs for model robustness. Zou et al. (2020) utilized reinforcement learning to generate adversarial examples, producing stable attacks with semantic-preserving adversarial examples. Cheng et al. (2020) proposed a novel adversarial augmentation method that minimizes the vicinal risk over virtual sentences sampled from a smoothly interpolated embedding space around the observed training sentence pairs. The adversarial data augmentation method substantially outperforms other data augmentation methods and achieves significant improvements in translation quality and robustness. For the better exploration of robust NMT, Michel and Neubig (2018) proposed an MTNT dataset, source sentences of which are collected from Reddit discussion, and contain several types of noise. Target referenced translations for each source sentence, in contrast, are clear from noise. Experiments showed that current NMT models perform badly on the MTNT dataset. As a result, this dataset can serve as a testbed for NMT robustness analysis.

## 3. Resources

### 3.1. Parallel data

Bilingual parallel corpora are the most important resources for NMT.

There are several publicly available corpora, such as the datasets provided by WMT,<sup>1</sup> IWSLT,<sup>2</sup> and WAT.<sup>3</sup> Table 3 lists the available domains and language pairs in these workshops.

Besides the aforementioned machine translation workshops, we also recommend OPUS<sup>4</sup> to search resources for training NMT models, which gathers parallel data for a large number of language pairs. We list the number of sentence pairs that are available for major languages to English in Table 4. OPUS also provides the OPUS-100 corpus for multilingual machine translation research (Zhang et al., 2020b), which is an English-centric multilingual corpus covering over 100 languages.

### 3.2. Monolingual data

Monolingual data are also valuable resources for NMT. The Common Crawl Foundation<sup>5</sup> provides open access to high quality crawled data for over 40 languages. The CCNet toolkit<sup>6</sup> (Wenzek et al., 2020) can be used to download and clean Common Crawl texts. Wikipedia provides database dump<sup>7</sup> that can be used to extract monolingual data, which can be download using WIKIEXTRACTOR.<sup>8</sup> WMT 2020 also provides several monolingual training data, which consists of data collected from NewsCrawl, NewsDiscussions, Europarl, NewsCommentary, CommonCrawl, and WikiDumps.

## 4. Tools

With the rapid advances of deep learning, many open-source deep learning frameworks have emerged, with TensorFlow (Abadi et al., 2016) and PyTorch (Paszke et al., 2019) as representative examples. At the same time, we have also witnessed the rapid development of open-source NMT toolkits, which significantly boosted the research progress of NMT. In this section, we will give a summarization of popular open-source NMT toolkits. Besides, we also introduce tools that are useful for evaluation, analysis, and data pre-processing.

### 4.1. Open-source NMT toolkits

We summarize some popular open-source NMT toolkits on GitHub in Table 5. The users can get the source codes of these toolkits directly from GitHub. We shall give a brief description of these projects.

**Tensor2Tensor.** TENSOR2TENSOR (Vaswani et al., 2018) is a library of deep learning models and datasets based on TensorFlow (Abadi et al., 2016). The library was mainly developed by the Google Brain team. TENSOR2TENSOR provides implementation of several NMT architectures (e.g., Transformer) for the translation task. The users can run TENSOR2TENSOR easily on CPU, GPU, and TPU, either locally or on Cloud.

**FairSeq.** FAIRSEQ (Ott et al., 2019) is a sequence modeling toolkit developed by Facebook AI Research. The toolkit is based on Pytorch (Paszke et al., 2019) and allows the users to train custom models for the translation task. FAIRSEQ implements traditional RNN-based models and Transformer models. Besides, it also includes CNN-based translation models (e.g., LightConv and DynamicConv).

**Nmt.** NMT (Luong et al., 2017) is a toolkit developed by Google Research. The toolkit implements the GNMT architecture (Wu et al., 2016). Besides, the NMT project also provides a nice tutorial for building a competitive NMT model from scratch. The codebase of NMT is high-quality and lightweight, which is friendly for users to add customized models.

<sup>2</sup> <http://iwslt.org/doku.php>.

<sup>3</sup> <http://lotus.kuee.kyoto-u.ac.jp/WAT/WAT2020/index.html>.

<sup>4</sup> <http://opus.nlpl.eu>.

<sup>5</sup> <https://commoncrawl.org/>.

<sup>6</sup> [https://github.com/facebookresearch/cc\\_net](https://github.com/facebookresearch/cc_net).

<sup>7</sup> <https://dumps.wikimedia.org>.

<sup>8</sup> <https://github.com/attardi/wikiextractor>.

<sup>1</sup> <http://www.statmt.org/wmt20/index.html>.



**OpenNMT.** OPENNMT is an open-source NMT toolkit developed by the collaboration of Harvard University and SYSTRAN. The toolkit currently maintains two implementations: OPENNMT-PY and OPENNMT-TF. OPENNMT is proven to be research-friendly and production-ready. The OpenNMT project also provides CTRANSLATE2 as a fast inference engine that supports both CPU and GPU.

**Sockeye.** SOCKEYE (Hieber et al., 2017) is a versatile sequence-to-sequence toolkit that is based on MXNet (Chen et al., 2016). SOCKEYE is maintained by Amazon and powers machine translation services such as Amazon Translate. The toolkit features state-of-the-art machine translation models and fast CPU inference, which is useful for both research and production.

**Nematus.** NEMATUS is an NMT toolkit developed by the NLP Group at the University of Edinburgh. The toolkit is based on TensorFlow and supports RNN-based NMT architectures as well as the TRANSFORMER architecture. In addition to the toolkits, NEMATUS also released high-performing NMT models covering 13 translation directions.

**Marian.** MARIAN (Junczys-Dowmunt et al., 2018) is an efficient and self-contained NMT framework currently being developed by the Microsoft Translator team. The framework is written entirely in C++ with minimal dependencies. Marian is widely deployed by many companies and organizations. For example, Microsoft Translator currently adopts Marian as its neural machine translation engine.

**THUMT.** THUMT (Zhang et al., 2017c) is an open-source toolkit for neural machine translation developed by the NLP Group at Tsinghua University. The toolkit includes Theano (Team et al., 2016), TensorFlow, and Pytorch implementations. It supports vanilla RNN-based and Transformer models and is easy for users to build new models. Furthermore, THUMT provides visualization analysis using layer-wise relevance propagation (Ding et al., 2017).

**NMT-Keras.** NMT-KERAS (Peris and Casacuberta, 2018) is a flexible toolkit for neural machine translation developed by the Pattern Recognition and Human Language Technology Research Center at Polytechnic University of Valencia. The toolkit is based on Keras which uses Theano or TensorFlow as the backend. NMT-KERAS emphasizes the development of advanced applications for NMT systems, such as interactive NMT and online learning. It also has been extended to other tasks including image and video captioning, sentence classification, and visual question answering.

**Neural Monkey** NEURAL MONKEY is an open-source neural machine translation and general sequence-to-sequence learning system. The toolkit is built on the TensorFlow library and provides a high-level API tailored for fast prototyping of complex architectures.

#### 4.2. Tools for evaluation and analysis

Manual evaluation of MT outputs is not only expensive but also impractical to scaling for large language pairs. On the contrary, automatic MT evaluation is inexpensive and language-independent, with BLEU (Papineni et al., 2002) as the representative automatic evaluation metric. Besides evaluation, there is also a need for analyzing MT outputs. We recommend the following tools for evaluating and analyzing MT output.

**SacreBLEU.** SACREBLEU<sup>9</sup> (Post, 2018) is a toolkit to compute shareable, comparable, and reproducible BLEU scores. SACREBLEU computes BLEU scores on detokenized outputs, using WMT standard tokenization. As a result, the scores are not affected by different processing tools. Besides, it can produce a short version string that facilitates cross-paper comparisons.

**COMPARE-MT.** COMPARE-MT<sup>10</sup> (Neubig et al., 2019) is a program to compare the outputs of multiple systems for language generation. In order to provide high-level analysis of outputs, it enables analysis of

accuracy of generation of particular types of words, bucketed histograms of sentence accuracies or counts based on salient characteristics, and so on.

**MT-COMPAReVAL.** MT-COMPAReVAL<sup>11</sup> is also a tool for comparison and evaluation of machine translations. It allows users to compare translations according to automatic metrics or quality comparison from the aspects of n-grams.

#### 4.3. Other tools

Asides from the above mentioned tools, we found the following toolkits are very useful for NMT research and deployment.

**MOSES.** MOSES<sup>12</sup> (Koehn et al., 2007) is a self-contained statistical machine translation toolkit. Besides SMT-related components, MOSES provides a large number of tools to clean and pre-process texts, which are also useful for training NMT models. MOSES also contains several easy-to-use scripts to analyze and evaluate MT outputs.

**SUBWORD-NMT.** SUBWORD-NMT<sup>13</sup> is an open-source toolkit for unsupervised word segmentation for neural machine translation and text generation. It adopts the Byte-Pair Encoding (BPE) algorithm proposed by (Sennrich et al., 2016c) and BPE dropout proposed by (Provilkov et al., 2019). It is the most commonly used toolkit to alleviate the out-of-vocabulary problem in NMT.

**SENTENCEPIECE.** SENTENCEPIECE<sup>14</sup> is a powerful unsupervised text segmentation toolkit. SENTENCEPIECE is written in C++ and provides APIs for other languages such as Python. SENTENCEPIECE implements the BPE algorithm (Sennrich et al., 2016c) and unigram language model (Kudo, 2018). Unlike SUBWORD-NMT, SENTENCEPIECE can learn to segment raw texts without additional pre-processing. As a result, SENTENCEPIECE is a suitable choice to segment multilingual texts.

## 5. Conclusion

Neural machine translation has become the dominant approach to machine translation in both research and practice. This article reviewed the widely used methods in NMT, including modeling, decoding, data augmentation, interpretation, as well as evaluation. We then summarize the resources and tools that are useful for NMT research.

Despite the great success achieved by NMT, there are still many problems to be explored. We list some important and challenging problems for NMT as follows:

- **Understanding NMT.** Although there are many attempts to analyze and interpret NMT, our understandings about NMT are still limited. Understanding how and why NMT produces its translation result is important to figure out the bottleneck and weakness of NMT models.
- **Designing better architectures.** Designing a new architecture that better than Transformer is beneficial for both NMT research and production. Furthermore, designing a new architecture that balances translation performance and computational complexity is also important.
- **Making full use of monolingual data.** Monolingual data are valuable resources. Despite the remarkable progress, we believe that there is still much room for NMT to make use of abundant monolingual data.
- **Prior knowledge integration.** Incorporating human knowledge into NMT is also an important problem. Although there is some progress, the results are far from satisfactory. How to convert discrete and continuous representations into each other is a problem of NMT that needs further exploration.

<sup>11</sup> <https://github.com/ondrejklejch/MT-CompareEval>.

<sup>12</sup> <https://github.com/moses-smt/mosesdecoder>.

<sup>13</sup> <https://github.com/rsennrich/subword-nmt>.

<sup>14</sup> <https://github.com/google/sentencepiece>.

<sup>9</sup> <https://github.com/mjpost/sacrebleu>.

<sup>10</sup> <https://github.com/neulab/compare-mt>.

## Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgements

This work was supported by the National Key R&D Program of China (No. 2017YFB0 202204), National Natural Science Foundation of China (No. 61925601, No. 61761166 008, No. 61772302), Beijing Academy of Artificial Intelligence, Huawei Noah's Ark Lab, and the NExT++ project supported by the National Research Foundation, Prime Ministers Office, Singapore under its IRC@Singapore Funding Initiative.

## References

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al., 2016. Tensorflow: a system for large-scale machine learning. In: 12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16), pp. 265–283.
- Aharoni, R., Goldberg, Y., 2017. Towards string-to-tree neural machine translation. In: Proceedings of ACL, pp. 132–140.
- Akoury, N., Krishna, K., Iyer, M., 2019. Syntactically supervised transformers for faster neural machine translation. In: Proceedings of ACL, pp. 1269–1281.
- Artetxe, M., Labaka, G., Agirre, E., 2017a. Learning bilingual word embeddings with (almost) no bilingual data. In: Proceedings of ACL, pp. 451–462.
- Artetxe, M., Labaka, G., Agirre, E., Cho, K., 2017b. Unsupervised Neural Machine Translation arXiv preprint arXiv:1710.11041.
- Artetxe, M., Labaka, G., Agirre, E., 2018. Unsupervised Statistical Machine Translation arXiv preprint arXiv:1809.01272.
- Artetxe, M., Labaka, G., Agirre, E., 2019. An Effective Approach to Unsupervised Machine Translation arXiv preprint arXiv:1902.01313.
- Ba, J.L., Kiros, J.R., Hinton, G.E., 2016. Layer Normalization arXiv preprint arXiv:1607.06450.
- Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K.-R., Samek, W., 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PLoS One 10, e0130140.
- Bahdanau, D., Cho, K., Bengio, Y., 2015. Neural machine translation by jointly learning to align and translate. In: Proceedings of ICLR.
- Baniata, L.H., Park, S., Park, S.-B., 2018. A multitask-based neural machine translation model with part-of-speech tags integration for Arabic dialects. Appl. Sci. 8, 2502.
- Bastings, J., Titov, I., Aziz, W., Marcheggiani, D., Sima'an, K., 2017. Graph convolutional encoders for syntax-aware neural machine translation. In: Proceedings of EMNLP, pp. 1957–1967.
- Bau, A., Belinkov, Y., Sajjad, H., Durrani, N., Dalvi, F., Glass, J., 2019. Identifying and controlling important neurons in neural machine translation. Proc. ICLR.
- Belinkov, Y., Bisk, Y., 2018. Synthetic and natural noise both break neural machine translation. In: Proceedings of ICLR.
- Bengio, Y., Simard, P., Frasconi, P., 1994. Learning long-term dependencies with gradient descent is difficult. IEEE Trans. Neural Network. 5, 157–166.
- Brown, P.F., Cocke, J., Della Pietra, S.A., Della Pietra, V.J., Jelinek, F., Lafferty, J., Mercer, R.L., Roossin, P.S., 1990. A statistical approach to machine translation. Comput. Ling. 16, 79–85.
- Bugliarello, E., Okazaki, N., 2020. Enhancing machine translation with dependency-aware self-attention. In: Proceedings of ACL, pp. 1618–1627.
- Caswell, I., Chelba, C., Grangier, D., 2019. Tagged back-translation. WMT 2019, 53.
- Chen, T., Li, M., Li, Y., Lin, M., Wang, N., Wang, M., Xiao, T., Xu, B., Zhang, C., Zhang, Z., 2016. Mxnet: a flexible and efficient machine learning library for heterogeneous distributed systems. In: Proceedings of NeurIPS, Workshop.
- Chen, H., Huang, S., Chiang, D., Dai, X., Chen, J., 2018. Combining character and word information in neural machine translation using a multi-level attention. In: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, vol. 1, pp. 1284–1293 (Long Papers).
- Chen, L., Zhang, Y., Zhang, R., Tao, C., Gan, Z., Zhang, H., Li, B., Shen, D., Chen, C., Carin, L., 2019. Improving sequence-to-sequence learning via optimal transport. In: Proceedings of ICLR.
- Chen, K., Wang, R., Utiyama, M., Sumita, E., 2020. Content word aware neural machine translation. In: Proceedings of ACL, pp. 358–364.
- Cheng, Y., Xu, W., He, Z., He, W., Wu, H., Sun, M., Liu, Y., 2016. Semi-supervised learning for neural machine translation. In: Proceedings of ACL, pp. 1965–1974.
- Cheng, Y., Tu, Z., Meng, F., Zhai, J., Liu, Y., 2018. Towards robust neural machine translation. In: Proceedings of ACL, pp. 1756–1766.
- Cheng, Y., Jiang, L., Macherey, W., 2019. Robust neural machine translation with doubly adversarial inputs. In: Proceedings of ACL, pp. 4324–4333.
- Cheng, Y., Jiang, L., Macherey, W., Eisenstein, J., AdvAug, 2020. Robust adversarial augmentation for neural machine translation. In: Proceedings of ACL, pp. 5961–5970.
- Cherry, C., Foster, G., Bapna, A., Firat, O., Macherey, W., 2018. Revisiting Character-Based Neural Machine Translation with Capacity and Compression arXiv preprint arXiv:1808.09943.
- Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y., 2014a. Learning phrase representations using rnn encoder–decoder for statistical machine translation. In: Proceedings of EMNLP, pp. 1724–1734.
- Cho, K., Van Merriënboer, B., Bahdanau, D., Bengio, Y., 2014b. On the Properties of Neural Machine Translation: Encoder-Decoder Approaches arXiv preprint arXiv:1409.1259.
- Choshen, L., Fox, L., Aizenbud, Z., Abend, O., 2020. On the weaknesses of reinforcement learning for neural machine translation. In: Proceedings of ICLR.
- Chung, J., Cho, K., Bengio, Y., 2016. A Character-Level Decoder without Explicit Segmentation for Neural Machine Translation arXiv preprint arXiv:1603.06147.
- Clark, K., Khandelwal, U., Levy, O., Manning, C.D., 2019. What does bert look at? an analysis of bert's attention. In: Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP, pp. 276–286.
- Conneau, A., Lample, G., Ranzato, M., Denoyer, L., Jégou, H., 2017. Word Translation without Parallel Data arXiv preprint arXiv:1710.04087.
- Devlin, J., Chang, M.-W., Lee, K., Toutanova, K., 2019. Bert: pre-training of deep bidirectional transformers for language understanding. In: Proceedings of NAACL-HLT, pp. 4171–4186.
- Ding, Y., Liu, Y., Luan, H., Sun, M., 2017. Visualizing and understanding neural machine translation. In: Proceedings of ACL, pp. 1150–1159.
- Ebrahimi, J., Lowd, D., Dou, D., 2018. On adversarial examples for character-level neural machine translation. In: Proceedings of COLING, pp. 653–663.
- Eduonov, S., Ott, M., Auli, M., Grangier, D., Ranzato, M., 2018. Classical structured prediction losses for sequence to sequence learning. In: Proceedings of NAACL-HLT, pp. 355–364.
- Eduonov, S., Ott, M., Auli, M., Grangier, D., 2018. Understanding Back-Translation at Scale arXiv preprint arXiv:1808.09381.
- Eduonov, S., Baevski, A., Auli, M., 2019. Pre-trained Language Model Representations for Language Generation arXiv preprint arXiv:1903.09722.
- Eriguchi, A., Hashimoto, K., Tsuruoka, Y., 2016. Tree-to-sequence attentional neural machine translation. In: Proceedings of ACL, pp. 823–833.
- Eriguchi, A., Tsuruoka, Y., Cho, K., 2017. Learning to parse and translate improves neural machine translation. In: Proceedings of ACL, pp. 72–78.
- Gao, Y., Nikolov, N.I., Hu, Y., Hahnloser, R.H., 2020. Character-level Translation with Self-Attention arXiv preprint arXiv:2004.14788.
- Gehring, J., Auli, M., Grangier, D., Yarats, D., Dauphin, Y.N., 2017. Convolutional Sequence to Sequence Learning arXiv preprint arXiv:1705.03122.
- Ghazvininejad, M., Levy, O., Liu, Y., Zettlemoyer, L., 2019. Mask-predict: parallel decoding of conditional masked language models. In: Proceedings of EMNLP-IJCNLP, pp. 6114–6123.
- Graves, A., Wayne, G., Danihelka, I., 2014. Neural Turing Machines arXiv preprint arXiv:1410.5401.
- Gu, J., Bradbury, J., Xiong, C., Li, V.O., Socher, R., 2018. Non-autoregressive neural machine translation. In: Proceedings of ICLR.
- Gu, J., Liu, Q., Cho, K., 2019. Insertion-based decoding with automatically inferred generation order. TACL 7, 661–676.
- Gulcehre, C., Firat, O., Xu, K., Cho, K., Bengio, Y., 2017. On integrating a language model into neural machine translation. Comput. Speech Lang 45, 137–148.
- Guo, J., Tan, X., He, D., Qin, T., Xu, L., Liu, T.-Y., 2019. Non-autoregressive neural machine translation with enhanced decoder input. In: Proceedings of AAAI, vol. 33, pp. 3723–3730.
- Gü, J., Shavaran, H.S., Sarkar, A., 2018. Top-down tree structured decoding with syntactic connections for neural machine translation and parsing. In: Proceedings of EMNLP, pp. 401–413.
- Hao, J., Wang, X., Shi, S., Zhang, J., Tu, Z., 2019. Multi-granularity self-attention for neural machine translation. In: Proceedings of EMNLP-IJCNLP, pp. 886–896.
- Hassan, H., Aue, A., Chen, C., Chowdhury, V., Clark, J., Federmann, C., Huang, X., Junczys-Dowmunt, M., Lewis, W., Li, M., et al., 2018. Achieving Human Parity on Automatic Chinese to English News Translation arXiv preprint arXiv:1803.05567.
- He, K., Zhang, X., Ren, S., Sun, J., 2016a. Deep residual learning for image recognition. Proc. CVPR 770–778.
- He, D., Xia, Y., Qin, T., Wang, L., Yu, N., Liu, T.-Y., Ma, W.-Y., 2016b. Dual learning for machine translation. In: Advances in NeurIPS, pp. 820–828.
- He, S., Tu, Z., Wang, X., Wang, L., Lyu, M., Shi, S., 2019. Towards understanding neural machine translation with word importance. In: Proceedings of EMNLP-IJCNLP, pp. 953–962.
- Hieber, F., Domhan, T., Denkowski, M., Vilar, D., Sokolov, A., Clifton, A., Post, M., 2017. Sockeye: A Toolkit for Neural Machine Translation arXiv preprint arXiv:1712.05690.
- Hoang, C.D.V., Haffari, G., Cohn, T., 2017. Towards decoding as continuous optimisation in neural machine translation. In: Proceedings of EMNLP, pp. 146–156.
- Hoang, V.C.D., Koehn, P., Haffari, G., Cohn, T., 2018. Iterative back-translation for neural machine translation. In: Proceedings of the 2nd Workshop on Neural Machine Translation and Generation, pp. 18–24.
- Hochreiter, S., Schmidhuber, J., 1997. Long Short-Term Memory, Neural Computation.
- Imamura, K., Fujita, A., Sumita, E., 2018. 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, pp. 55–63.
- Junczys-Dowmunt, M., Dwojak, T., Hoang, H., 2016. Is Neural Machine Translation Ready for Deployment? a Case Study on 30 Translation Directions arXiv preprint arXiv:1610.01108.
- Junczys-Dowmunt, M., Grundkiewicz, R., Dwojak, T., Hoang, H., Heafield, K., Neckermann, T., Seide, F., Germann, U., Aji, A.F., Bogoychev, N., Martins, A.F.T., Birch, A., 2018. Marian: fast neural machine translation in C++. In: Proceedings of ACL System Demonstrations, pp. 116–121.
- Kaiser, L., Gomez, A.N., Chollet, F., 2017. Depthwise Separable Convolutions for Neural Machine Translation arXiv preprint arXiv:1706.03059.

- Kalchbrenner, N., Blunsom, P., 2013. Recurrent continuous translation models. In: *Proceedings of EMNLP*, pp. 1700–1709.
- Kalchbrenner, N., Espeholt, L., Simonyan, K., Oord, A.v. d., Graves, A., Kavukcuoglu, K., 2016. Neural Machine Translation in Linear Time arXiv preprint arXiv:1610.10099.
- Karakanta, A., Dehdari, J., van Genabith, J., 2018. Neural machine translation for low-resource languages without parallel corpora. *Mach. Translat.* 32, 167–189.
- Kim, Y., Rush, A.M., 2016. Sequence-level knowledge distillation. In: *Proceedings of EMNLP*, pp. 1317–1327.
- Kingma, D., Ba, J., Adam, 2014. A Method for Stochastic Optimization arXiv preprint arXiv:1412.6980.
- Koehn, P., Och, F.J., Marcu, D., 2003. Statistical Phrase-Based Translation, Technical Report. UNIVERSITY OF SOUTHERN CALIFORNIA MARINA DEL REY INFORMATION SCIENCES INST.
- Koehn, P., Hoang, H., Birch, A., Callison-Burch, C., Federico, M., Bertoldi, N., Cowan, B., Shen, W., Moran, C., Zens, R., et al., 2007. Moses: open source toolkit for statistical machine translation. In: *Proceedings of ACL on Interactive Poster and Demonstration Sessions*, pp. 177–180.
- Kudo, T., 2018. Subword Regularization: Improving Neural Network Translation Models with Multiple Subword Candidates arXiv preprint arXiv:1804.10959.
- Kumar, S., Tsvetkov, Y., 2019. Von mises-Fisher loss for training sequence to sequence models with continuous outputs. In: *Proceedings of ICLR*.
- Lample, G., Conneau, A., 2019. Cross-lingual Language Model Pretraining arXiv preprint arXiv:1901.07291.
- Lample, G., Conneau, A., Denoyer, L., Ranzato, M., 2017. Unsupervised Machine Translation Using Monolingual Corpora Only arXiv preprint arXiv:1711.00043.
- Lample, G., Ott, M., Conneau, A., Denoyer, L., Ranzato, M., 2018. Phrase-based & Neural Unsupervised Machine Translation arXiv preprint arXiv:1804.07755.
- Lee, J., Cho, K., Hofmann, T., 2017. Fully character-level neural machine translation without explicit segmentation. *Trans. Assoc. Comput. Ling.* 5, 365–378.
- Lee, J., Mansimov, E., Cho, K., 2018. Deterministic non-autoregressive neural sequence modeling by iterative refinement. In: *Proceedings of EMNLP*, pp. 1173–1182.
- Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V., Zettlemoyer, L., 2019. Bart: Denoising Sequence-To-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension arXiv preprint arXiv:1910.13461.
- Li, X., Liu, L., Tu, Z., Shi, S., Meng, M., 2018. Target foresight based attention for neural machine translation. In: *Proceedings of NAACL-HLT*, pp. 1380–1390.
- Libovický, J., Held, J., 2018. End-to-end non-autoregressive neural machine translation with connectionist temporal classification. In: *Proceedings of EMNLP*, pp. 3016–3021.
- Liu, Y., Liu, Q., Lin, S., 2006. Tree-to-string alignment template for statistical machine translation. In: *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Sydney, Australia, pp. 609–616. <https://doi.org/10.3115/1220175.1220252>. URL: <https://www.aclweb.org/anthology/P06-1077>.
- Liu, L., Utiyama, M., Finch, A., Sumita, E., 2016. Agreement on target-bidirectional neural machine translation. In: *Proceedings of NAACL-HLT*, pp. 411–416.
- Liu, H., Ma, M., Huang, L., Xiong, H., He, Z., 2019a. Robust neural machine translation with joint textual and phonetic embedding. In: *Proceedings of ACL*, pp. 3044–3049.
- Liu, X., Wong, D.F., Liu, Y., Chao, L.S., Xiao, T., Zhu, J., 2019b. Shared-private bilingual word embeddings for neural machine translation. In: *Proceedings of ACL*, pp. 3613–3622.
- Liu, Y., Gu, J., Goyal, N., Li, X., Edunov, S., Ghazvininejad, M., Lewis, M., Zettlemoyer, L., 2020. Multilingual Denoising Pre-training for Neural Machine Translation arXiv preprint arXiv:2001.08210.
- Luong, M.-T., Manning, C.D., 2016. Achieving Open Vocabulary Neural Machine Translation with Hybrid Word-Character Models arXiv preprint arXiv:1604.00788.
- Luong, M.-T., Pham, H., Manning, C.D., 2015. Effective Approaches to Attention-Based Neural Machine Translation arXiv preprint arXiv:1508.04025.
- Luong, M., Brevdo, E., Zhao, R., 2017. Neural machine translation (seq2seq) tutorial. <http://github.com/tensorflow/nmt>.
- Mehri, S., Sigal, L., 2018. Middle-out decoding. In: *Advances in NeurIPS*, pp. 5518–5529.
- Michel, P., Neubig, G., 2018. Mtn: a testbed for machine translation of noisy text. In: *Proceedings of EMNLP*, pp. 543–553.
- Morishita, M., Suzuki, J., Nagata, M., 2018. Improving neural machine translation by incorporating hierarchical subword features. In: *Proceedings of COLING*, pp. 618–629.
- Neubig, G., Dou, Z.-Y., Hu, J., Michel, P., Pruthi, D., Wang, X., 2019. compare-mt: a tool for holistic comparison of language generation systems. In: *Proceedings of NAACL-HLT (Demonstrations)*, pp. 35–41.
- Ott, M., Edunov, S., Baevski, A., Fan, A., Gross, S., Ng, N., Grangier, D., Auli, M., 2019. fairseq: a fast, extensible toolkit for sequence modeling. In: *Proceedings of NAACL-HLT (Demonstrations)*, pp. 48–53.
- Papineni, K., Roukos, S., Ward, T., Zhu, W., 2002. Bleu: a method for automatic evaluation of machine translation. In: *Proceedings of ACL*.
- Passban, P., Liu, Q., Way, A., 2018. Improving Character-Based Decoding Using Target-Side Morphological Information for Neural Machine Translation arXiv preprint arXiv:1804.06506.
- Paszkiewicz, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimeshain, N., Antiga, L., et al., 2019. Pytorch: an imperative style, high-performance deep learning library. In: *Advances in NeurIPS*, pp. 8026–8037.
- Peris, A., Casacuberta, F., 2018. Nmt-keras: a very flexible toolkit with a focus on interactive nmt and online learning. *Prague Bull. Math. Linguist.* 111, 113–124.
- Peters, M.E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., Zettlemoyer, L., 2018. Deep Contextualized Word Representations arXiv preprint arXiv:1802.05365.
- Poncelas, A., Shterionov, D., Way, A., de Buy Wenniger, G., Passban, P., 2018. Investigating Backtranslation in Neural Machine Translation arXiv preprint arXiv:1804.06189.
- Post, M., 2018. A Call for Clarity in Reporting Bleu Scores arXiv preprint arXiv:1804.08771.
- Provlkov, I., Emelianenko, D., Voita, E., 2019. Bpe-dropout: Simple and Effective Subword Regularization arXiv preprint arXiv:1910.13267.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., 2019. Language models are unsupervised multitask learners. *OpenAI Blog* 1, 9.
- Raganato, A., Tiedemann, J., 2018. An analysis of encoder representations in transformer-based machine translation. In: *Proceedings of EMNLP Workshop*, pp. 287–297.
- Ranzato, M., Chopra, S., Auli, M., Zaremba, W., 2015. Sequence Level Training with Recurrent Neural Networks arXiv preprint arXiv:1511.06732.
- Ranzato, M., Chopra, S., Auli, M., Zaremba, W., 2016. Sequence level training with recurrent neural networks. In: *Proceedings of ICLR*.
- Ren, S., Zhang, Z., Liu, S., Zhou, M., Ma, S., 2019. Unsupervised neural machine translation with smt as posterior regularization. In: *Proceedings of the AAAI*, vol. 33, pp. 241–248.
- Ren, S., Wu, Y., Liu, S., Zhou, M., Ma, S., 2020. A retrieve-and-rewrite initialization method for unsupervised machine translation. In: *Proceedings of ACL*, pp. 3498–3504.
- Ribeiro, M.T., Singh, S., Guestrin, C., 2018. Semantically equivalent adversarial rules for debugging nlp models. In: *Proceedings of ACL*, pp. 856–865.
- Sennrich, R., Haddow, B., 2016. Linguistic input features improve neural machine translation. In: *Proceedings of WMT*, pp. 83–91.
- Sennrich, R., Haddow, B., Birch, A., 2016a. Edinburgh neural machine translation systems for wmt 16. In: *Proceedings of WMT*, pp. 371–376.
- Sennrich, R., Haddow, B., Birch, A., 2016b. Improving neural machine translation models with monolingual data. In: *Proceedings of ACL*, pp. 86–96.
- Sennrich, R., Haddow, B., Birch, A., 2016c. Neural machine translation of rare words with subword units. In: *Proceedings of ACL*.
- Sennrich, R., Birch, A., Currey, A., Hermann, U., Haddow, B., Heafield, K., Barone, A.V.M., Williams, P., 2017. The university of edinburgh's neural mt systems for wmt17. In: *Proceedings of WMT*, pp. 389–399.
- Shao, C., Feng, Y., Zhang, J., Meng, F., Chen, X., Zhou, J., 2019. Retrieving sequential information for non-autoregressive neural machine translation. In: *Proceedings of ACL*, pp. 3013–3024.
- Shen, S., Cheng, Y., He, Z., He, W., Wu, H., Sun, M., Liu, Y., 2016. Minimum risk training for neural machine translation. In: *Proceedings of ACL*, pp. 1683–1692.
- Sieglmann, H.T., Sontag, E.D., 1995. On the computational power of neural nets. *J. Comput. Syst. Sci.* 50, 132–150.
- Song, K., Tan, X., Qin, T., Lu, J., Liu, T.-Y., 2019. Mass: Masked Sequence to Sequence Pre-training for Language Generation arXiv preprint arXiv:1905.02450.
- Stahlberg, F., Saunders, D., Byrne, B., 2018. An operation sequence model for explainable neural machine translation. In: *Proceedings of EMNLP Workshop*, pp. 175–186.
- Stern, M., Chan, W., Kiros, J., Uszkoreit, J., 2019. Insertion transformer: flexible sequence generation via insertion operations. In: *Proceedings of ICML*, pp. 5976–5985.
- Strobelt, H., Gehrmann, S., Behrisch, M., Perer, A., Pfister, H., Rush, A.M., 2019. Seq2seqvis: a visual debugging tool for sequence-to-sequence models. *IEEE Trans. Visual. Comput. Graph.* 25, 353–363.
- Su, J., Zhang, X., Lin, Q., Qin, Y., Yao, J., Liu, Y., 2019. Exploiting reverse target-side contexts for neural machine translation via asynchronous bidirectional decoding. *Artif. Intell.* 277, 103168.
- Sun, H., Wang, R., Chen, K., Utiyama, M., Sumita, E., Zhao, T., 2019. Unsupervised bilingual word embedding agreement for unsupervised neural machine translation. In: *Proceedings of ACL*, pp. 1235–1245.
- Sutskever, I., Vinyals, O., Le, Q.V., 2014. Sequence to sequence learning with neural networks. In: *Proceedings of NeurIPS*, pp. 3104–3112.
- Team, T.T.D., Al-Rfou, R., Alain, G., Almahairi, A., Angermueller, C., Bahdanau, D., Ballas, N., Bastien, F., Bayer, J., Belikov, A., et al., 2016. Theano: A python Framework for Fast Computation of Mathematical Expressions arXiv preprint arXiv:1605.02688.
- Tiedemann, J., 2016. Opus-parallel corpora for everyone. *Baltic J. Mod. Comput.* 384.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I., 2017. Attention is all you need. In: *Proceedings of NeurIPS*, pp. 5998–6008.
- Vaswani, A., Bengio, S., Brevdo, E., Chollet, F., Gomez, A., Gouws, S., Jones, L., Kaiser, Ł., Kalchbrenner, N., Parmar, N., Sepassi, R., Shazeer, N., Uszkoreit, J., 2018. Tensor2Tensor for neural machine translation. In: *Proceedings of AMTA*, pp. 193–199.
- Vincent, P., Larochelle, H., Bengio, Y., Manzagol, P.-A., 2008. Extracting and composing robust features with denoising autoencoders. In: *Proceedings of ICML*, pp. 1096–1103.
- Voita, E., Sennrich, R., Titov, I., 2019. The bottom-up evolution of representations in the transformer: a study with machine translation and language modeling objectives. In: *Proceedings of EMNLP-IJCNLP*, pp. 4396–4406.
- Wang, C., Zhang, J., Chen, H., 2018a. Semi-autoregressive neural machine translation. In: *Proceedings of EMNLP*, pp. 479–488.
- Wang, X., Pham, H., Yin, P., Neubig, G., 2018b. A tree-based decoder for neural machine translation. In: *Proceedings of EMNLP*, pp. 4772–4777.
- Wang, Y., Tian, F., He, D., Qin, T., Zhai, C., Liu, T.-Y., 2019a. Non-autoregressive machine translation with auxiliary regularization. In: *Proceedings of AAAI*, vol. 33, pp. 5377–5384.
- Wang, S., Liu, Y., Wang, C., Luan, H., Sun, M., 2019b. Improving Back-Translation with Uncertainty-Based Confidence Estimation arXiv preprint arXiv:1909.00157.



- Wang, C., Cho, K., Gu, J., 2020. Neural machine translation with byte-level subwords. In: Proceedings of the AAAI Conference on Artificial Intelligence, vol. 34, pp. 9154–9160.
- Wei, B., Wang, M., Zhou, H., Lin, J., Sun, X., 2019. Imitation learning for non-autoregressive neural machine translation. In: Proceedings of ACL, pp. 1304–1312.
- Weller-Di Marco, M., Fraser, A., 2020. Modeling word formation in English–German neural machine translation. In: Proceedings of ACL, pp. 4227–4232.
- Wenzek, G., Lachaux, M.-A., Conneau, A., Chaudhary, V., Guzmán, F., Joulin, A., Grave, É., Ccnet, 2020. Extracting high quality monolingual datasets from web crawl data. In: Proceedings of the 12th Language Resources and Evaluation Conference, pp. 4003–4012.
- Wieting, J., Berg-Kirkpatrick, T., Gimpel, K., Neubig, G., 2019. Beyond BLEU: training neural machine translation with semantic similarity. In: Proceedings of ACL, pp. 4344–4355.
- Wu, Y., Schuster, M., Chen, Z., Le, Q.V., et al., 2016. Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation arXiv preprint arXiv:1609.08144.
- Wu, S., Zhang, D., Yang, N., Li, M., Zhou, M., 2017. Sequence-to-dependency neural machine translation. In: Proceedings of ACL, pp. 698–707.
- Wu, L., Tian, F., Qin, T., Lai, J., Liu, T.-Y., 2018. A study of reinforcement learning for neural machine translation. In: Proceedings of EMNLP, pp. 3612–3621.
- Wu, F., Fan, A., Baevski, A., Dauphin, Y.N., Auli, M., 2019a. Pay Less Attention with Lightweight and Dynamic Convolutions arXiv preprint arXiv:1901.10430.
- Wu, J., Wang, X., Wang, W.Y., 2019b. Extract and Edit: an Alternative to Back-Translation for Unsupervised Neural Machine Translation arXiv preprint arXiv:1904.02331.
- Yang, Z., Chen, L., Le Nguyen, M., 2018a. Regularizing forward and backward decoding to improve neural machine translation. In: Proceedings of International Conference on Knowledge and Systems Engineering. KSE, pp. 73–78.
- Yang, Z., Chen, W., Wang, F., Xu, B., 2018b. Unsupervised Neural Machine Translation with Weight Sharing arXiv preprint arXiv:1804.09057.
- Yang, Z., Cheng, Y., Liu, Y., Sun, M., 2019a. Reducing word omission errors in neural machine translation: a contrastive learning approach. In: Proceedings of ACL, pp. 6191–6196.
- Yang, X., Liu, Y., Xie, D., Wang, X., Balasubramanian, N., 2019b. Latent part-of-speech sequences for neural machine translation. In: Proceedings of EMNLP-IJCNLP, pp. 780–790.
- Yang, J., Ma, S., Zhang, D., Li, Z., Zhou, M., 2020. Improving neural machine translation with soft template prediction. In: Proceedings of WMT, pp. 5979–5989.
- Yun, C., Bhojanapalli, S., Rawat, A.S., Reddi, S., Kumar, S., 2020. Are transformers universal approximators of sequence-to-sequence functions?. In: Proceedings of ICLR.
- Zhang, J., Zong, C., 2016. Exploiting source-side monolingual data in neural machine translation. In: Proceedings of EMNLP, pp. 1535–1545.
- Zhang, J., Zong, C., 2020. Neural Machine Translation: Challenges, Progress and Future arXiv preprint arXiv:2004.05809.
- Zhang, M., Liu, Y., Luan, H., Sun, M., 2017a. Adversarial training for unsupervised bilingual lexicon induction. In: Proceedings of ACL, pp. 1959–1970.
- Zhang, J., Liu, Y., Luan, H., Xu, J., Sun, M., 2017b. Prior knowledge integration for neural machine translation using posterior regularization. In: Proceedings of ACL, pp. 1514–1523.
- Zhang, J., Ding, Y., Shen, S., Cheng, Y., Sun, M., Luan, H., Liu, Y., 2017c. Thumt: an Open Source Toolkit for Neural Machine Translation arXiv preprint arXiv:1706.06415.
- Zhang, X., Su, J., Qin, Y., Liu, Y., Ji, R., Wang, H., 2018. Asynchronous bidirectional decoding for neural machine translation. In: Proceedings of AAAI, pp. 5698–5705.
- Zhang, B., Xiong, D., Su, J., Luo, J., 2019a. Future-aware knowledge distillation for neural machine translation. IEEE/ACM Trans. Audio Speech Lang. Proc. 27, 2278–2287.
- Zhang, Z., Wu, S., Liu, S., Li, M., Zhou, M., Xu, T., 2019b. Regularizing neural machine translation by target-bidirectional agreement. Proc. AAAI 33, 443–450.
- Zhang, J., Zhou, L., Zhao, Y., Zong, C., 2020a. Synchronous bidirectional inference for neural sequence generation. Artif. Intell. 281, 103234.
- Zhang, B., Williams, P., Titov, I., Sennrich, R., 2020b. Improving Massively Multilingual Neural Machine Translation and Zero-Shot Translation arXiv preprint arXiv:2004.11867.
- Zhao, Z., Dua, D., Singh, S., 2018. Generating natural adversarial examples. In: Proceedings of ICLR.
- Zheng, Z., Zhou, H., Huang, S., Mou, L., Dai, X., Chen, J., Tu, Z., 2018. Modeling past and future for neural machine translation. Trans. Assoc. Comput. Ling. 6, 145–157.
- Zheng, Z., Huang, S., Tu, Z., Dai, X.-Y., Jiajun, C., 2019. Dynamic past and future for neural machine translation. In: Proceedings of EMNLP-IJCNLP, pp. 930–940.
- Zheng, Z., Zhou, H., Huang, S., Li, L., Dai, X.-Y., Chen, J., 2020. Mirror-generative neural machine translation. In: Proceedings of ICLR.
- Zhou, J., Xu, W., 2015. End-to-end learning of semantic role labeling using recurrent neural networks. In: Proceedings of ACL, pp. 1127–1137.
- Zhou, L., Zhang, J., Zong, C., 2019a. Synchronous Bidirectional Neural Machine Translation. TACL.
- Zhou, L., Zhang, J., Zong, C., Yu, H., 2019b. Sequence generation: from both sides to the middle. In: Proceedings of IJCAI, pp. 5471–5477.
- Zhou, C., Gu, J., Neubig, G., 2019c. Understanding knowledge distillation in non-autoregressive machine translation. In: Proceedings of ICLR.
- Zhu, J., Xia, Y., Wu, L., He, D., Qin, T., Zhou, W., Li, H., Liu, T.-Y., 2020. Incorporating Bert into Neural Machine Translation arXiv preprint arXiv:2002.06823.
- Zoph, B., Yuret, D., May, J., Knight, K., 2016. Transfer Learning for Low-Resource Neural Machine Translation arXiv preprint arXiv:1604.02201.
- Zou, W., Huang, S., Xie, J., Dai, X., Chen, J., 2020. A reinforced generation of adversarial examples for neural machine translation. In: Proceedings of ACL, pp. 3486–3497.