

# Neural Abstractive Text Summarization with Sequence-to-Sequence Models

Tian Shi, Yaser Keneshloo, Naren Ramakrishnan, Chandan K. Reddy, *Senior Member, IEEE*

**Abstract**—In the past few years, neural abstractive text summarization with sequence-to-sequence (seq2seq) models have gained a lot of popularity. Many interesting techniques have been proposed to improve the seq2seq models, making them capable of handling different challenges, such as saliency, fluency and human readability, and generate high-quality summaries. Generally speaking, most of these techniques differ in one of these three categories: network structure, parameter inference, and decoding/generation. There are also other concerns, such as efficiency and parallelism for training a model. In this paper, we provide a comprehensive literature and technical survey on different seq2seq models for abstractive text summarization from viewpoint of network structures, training strategies, and summary generation algorithms. Many models were first proposed for language modeling and generation tasks, such as machine translation, and later applied to abstractive text summarization. Therefore, we also provide a brief review of these models. As part of this survey, we also develop an open source library, namely Neural Abstractive Text Summarizer (NATS) toolkit, for the abstractive text summarization. An extensive set of experiments have been conducted on the widely used CNN/Daily Mail dataset to examine the effectiveness of several different neural network components. Finally, we benchmark two models implemented in NATS on two recently released datasets, i.e., Newsroom and Bytcup.

**Index Terms**—Abstractive text summarization, sequence-to-sequence models, attention model, pointer-generator network, deep reinforcement learning, beam search.

## I. INTRODUCTION

In the modern era of big data, retrieving useful information from a large number of textual documents is a challenging task, due to the unprecedented growth in the availability of blogs, news articles, and reports are explosive. Automatic text summarization provides an effective solution for summarizing these documents. The task of the text summarization is to condense long documents into short summaries while preserving the important information and meaning of the documents [1, 2]. Having the short summaries, the text content can be retrieved, processed and digested effectively and efficiently.

Generally speaking, there are two ways to do text summarization: Extractive and Abstractive [3]. A method is considered to be *extractive* if words, phrases, and sentences in the summaries are selected from the source articles [4, 5, 6, 2, 7, 8, 9, 10]. They are relatively simple and can produce grammatically correct sentences. The generated summaries usually persist salient information of source articles and have a

good matching with human-written summaries [5, 11, 12, 13]. On the other hand, abstractive text summarization has attracted many attentions since it is capable of generating novel words using language generation models grounded on representations of source documents [14, 15]. Thus, they have a strong potential of producing high-quality summaries that are verbally innovative and can also easily incorporate external knowledge [12]. In this category, many deep neural network based models have achieved better performance in terms of the commonly used evaluation measures (such as *ROUGE* [16] score) compared to traditional extractive approaches [17, 18]. In this paper, we primarily focus on the recent advances of sequence-to-sequence (seq2seq) models for the task of abstractive text summarization.

### A. Seq2seq Models and Pointer-Generator Network

Seq2seq models (see Fig. 2) [19, 20] have been successfully applied to a variety of natural language processing (NLP) tasks, such as machine translation [21, 22, 23, 24, 25], headline generation [15, 26, 27], text summarization [12, 14], and speech recognition [28, 29, 30]. Inspired by the success of neural machine translation (NMT) [23], Rush *et al.* [15] first introduced a neural attention seq2seq model with an attention based encoder and a neural network language model (NNLM) decoder to the abstractive sentence summarization task, which has achieved a significant performance improvement over conventional methods. Chopra *et al.* [26] further extended this model by replacing the feed-forward NNLM with a recurrent neural network (RNN). The model is also equipped with a convolutional attention-based encoder and a RNN (Elman [31] or LSTM [32]) decoder, and outperforms other state-of-the-art models on a commonly used benchmark dataset, i.e., the Gigaword corpus. Nallapati *et al.* [14] introduced several novel elements to the RNN encoder-decoder architecture to address critical problems in the abstractive text summarization, including using the following (i) feature-rich encoder to capture keywords, (ii) a switching generator-pointer to model out-of-vocabulary (OOV) words, and (iii) the hierarchical attention to capture hierarchical document structures. They also established benchmarks for these models on a CNN/Daily Mail dataset [33, 34], which consists of pairs of news articles and multi-sentence highlights (summaries). Before this dataset was introduced, many abstractive text summarization models have concentrated on compressing short documents to single sentence summaries [15, 26]. For the task of summarizing long documents into multi-sentence summaries, these models have several shortcomings: 1) They cannot accurately reproduce

T. Shi, Y. Keneshloo, N. Ramakrishnan, and C. K. Reddy are with the Discovery Analytics Center, Department of Computer Science at Virginia Tech, Arlington, VA. {tshi,yaserkl}@vt.edu, {naren,reddy}@cs.vt.edu. Corresponding author: tshi@vt.edu.

the salient information of source documents. 2) They cannot efficiently handle OOV words. 3) They tend to suffer from word- and sentence-level repetitions and generating unnatural summaries. To tackle the first two challenges, See *et al.* [12] proposed a pointer-generator network that implicitly combines the abstraction with the extraction. This pointer-generator architecture can copy words from source texts via a pointer and generate novel words from a vocabulary via a generator. With the pointing/copying mechanism [35, 36, 37, 38, 39, 40, 41], factual information can be reproduced accurately and OOV words can also be taken care in the summaries. Many subsequent studies that achieved state-of-the-art performance have also demonstrated the effectiveness of the pointing/copying mechanism [17, 18, 42, 43]. The third problem has been addressed by the coverage mechanism [12], intra-temporal and intra-decoder attention mechanisms [17], and some other heuristic approaches, like forcing a decoder to never output the same trigram more than once during testing [17].

### B. Training Strategies

There are two other non-trivial issues with the current seq2seq framework, i.e., *exposure bias* and *inconsistency of training and testing measurements* [44, 45, 46, 47]. Based on the neural probabilistic language model [48], seq2seq models are usually trained by maximizing the likelihood of ground-truth tokens given their previous ground-truth tokens and hidden states (Teacher Forcing algorithm [44, 49], see Fig. 9). However, at testing time (see Fig. 8), previous ground-truth tokens are unknown, and they are replaced with tokens generated by the model itself. Since the generated tokens have never been exposed to the decoder during training, the decoding error can accumulate quickly during the sequence generation. This is known as *exposure bias* [46]. The other issue is the *mismatch of measurements*. Performance of seq2seq models is usually estimated with non-differentiable evaluation metrics, such as ROUGE [16] and BLEU [50] scores, which are inconsistent with the log-likelihood function (cross-entropy loss) used in the training phase. These problems are alleviated by the curriculum learning and reinforcement learning (RL) approaches.

1) *Training with Curriculum and Reinforcement Learning Approaches*: Bengio *et al.* [44] proposed a curriculum learning approach, known as *scheduled sampling*, to slowly change the input of the decoder from ground-truth tokens to model generated ones. Thus, the proposed meta-algorithm bridges the gap between training and testing. It is a practical solution for avoiding the exposure bias. Ranzato *et al.* [46] proposed a sequence level training algorithm, called MIXER (Mixed Incremental Cross-Entropy Reinforce), which consists of the cross entropy training, REINFORCE [51] and curriculum learning [44]. REINFORCE algorithm can make use of any user-defined task specific reward (e.g., non-differentiable evaluation metrics), therefore, combining with curriculum learning, the proposed model is capable of addressing both issues of seq2seq models. However, REINFORCE suffers from the high variance of gradient estimators and instability during training [52, 53, 22]. Bahdanau *et al.* [52] proposed an actor-critic based RL method

which has relatively lower variance for gradient estimators. In the actor-critic method, an additional critic network is trained to compute value functions given the policy from the actor network (a seq2seq model), and the actor network is trained based on the estimated value functions (assumed to be exact) from the critic network. On the other hand, Rennie *et al.* [53] introduced a self-critical sequence training method (SCST) which has a lower variance compared to the REINFORCE algorithm and does not need the second critic network.

2) *Applications to Abstractive Text Summarization*: RL algorithms for training seq2seq models have achieved success in a variety of language generation tasks, such as image captioning [53], machine translation [52], and dialogue generation [54]. Specific to the abstractive text summarization, Lin *et al.* [55] introduced a coarse-to-fine attention framework for the purpose of summarizing long documents. Their model parameters were learned with REINFORCE algorithm. Zhang *et al.* [56] used REINFORCE algorithm and the curriculum learning strategy for the sentence simplification task. Paulus *et al.* [17] first applied the self-critic policy gradient algorithm to training their seq2seq model with the copying mechanism and obtained the state-of-the-art performance in terms of ROUGE scores [16]. They proposed a mixed objective function that combines the RL loss with the traditional cross-entropy loss. Thus, their method can both leverage the non-differentiable evaluation metrics and improve the readability. Celikyilmaz *et al.* [18] introduced a novel deep communicating agents method for abstractive summarization, where they also adopted the RL loss in their objective function. Pasunuru *et al.* [57] applied the self-critic policy gradient algorithm to train the pointer-generator network. They also introduced two novel rewards (i.e., saliency and entailment rewards) in addition to ROUGE metric to keep the generated summaries salient and logically entailed. Li *et al.* [58] proposed a training framework based on the actor-critic method, where the actor network is an attention-based seq2seq model, and the critic network consists of a maximum likelihood estimator and a global summary quality estimator that is used to distinguish the generated and ground-truth summaries via a neural network binary classifier. Chen *et al.* [59] proposed a compression-paraphrase multi-step procedure, for abstractive text summarization, which first extracts salient sentences from documents and then rewrites them. In their model, they used an advantage actor-critic algorithm to optimize the sentence extractor for a better extraction strategy. Keneshloo *et al.* [47] conducted a comprehensive summary of various RL methods and their applications in training seq2seq models for different NLP tasks. They also implemented these RL algorithms in an open source library<sup>1</sup> constructed using the pointer-generator network [12] as the base model.

### C. Beyond RNN

Most of the prevalent seq2seq models that have attained state-of-the-art performance for sequence modeling and language generation tasks are RNN, especially long short-term memory (LSTM) [32] and gated recurrent unit (GRU) [60],

<sup>1</sup><https://github.com/yaserkl/RLSeq2Seq/>

based encoder-decoder models [19, 23]. Standard RNN models are difficult to train due to the vanishing and exploding gradients problems [61]. LSTM is a solution for vanishing gradients problem, but still does not address the exploding gradients issue. This issue is recently solved using a gradient norm clipping strategy [62]. Another critical problem of RNN based models is the computation constraint for long sequences due to their inherent sequential dependence nature. In other words, the current hidden state in a RNN is a function of previous hidden states. Because of such dependence, RNN cannot be parallelized within a sequence along the time-step dimension (see Fig. 2) during training and evaluation, and hence training them becomes major challenge for long sequences due to the computation time and memory constraints of GPUs [63].

Recently, it has been found that the convolutional neural network (CNN) [64] based encoder-decoder models have the potential to alleviate the aforementioned problem, since they have better performance in terms of the following three considerations [65, 66, 67]. 1) A model can be parallelized during training and evaluation. 2) The computational complexity of the model is linear with respect to the length of sequences. 3) The model has short paths between pairs of input and output tokens, so that it can propagate gradient signals more efficiently [68]. Kalchbrenner *et al.* [65] introduced a *ByteNet* model which adopts the one-dimensional convolutional neural network of fixed depth to both the encoder and the decoder [69]. The decoder CNN is stacked on top of the hidden representation of the encoder CNN, which ensures a shorter path between input and output. The proposed ByteNet model has achieved state-of-the-art performance on a character-level machine translation task with parallelism and linear-time computational complexity [65]. Bradbury *et al.* [66] proposed a quasi-recurrent neural network (QRNN) encoder-decoder architecture, where both encoder and decoder are composed of convolutional layers and so-called ‘dynamic average pooling’ layers [70, 66]. The convolutional layers allow computations to be completely parallel across both mini-batches and sequence time-step dimensions, while they require less amount of time compared with computation demands for LSTM despite the sequential dependence still presents in the pooling layers [66]. This framework has demonstrated to be effective by outperforming LSTM-based models on a character-level machine translation task with a significantly higher computational speed. Recently, Gehring *et al.* [71, 67, 72] attempted to build CNN based seq2seq models and apply them to large-scale benchmark datasets for sequence modeling. In [71], the authors proposed a convolutional encoder model, in which the encoder is composed of a succession of convolutional layers, and demonstrated its strong performance for machine translation. They further constructed a convolutional seq2seq architecture by replacing the LSTM decoder with a CNN decoder and bringing in several novel elements, including gated linear units [73] and multi-step attention [67]. The model also enables computations of all network elements parallelized, thus training and decoding can be much faster than the RNN models. It also achieved state-of-the-art performance on several machine translation benchmark datasets. Vaswani *et al.* [63] further constructed a novel network ar-

chitecture, namely, Transformer, which only depends on feed-forward networks and the attention mechanism. It has achieved state-of-the-art performance in machine translation task with significantly less training time. Currently, ConvS2S model [72] has been applied to the abstractive document summarization and outperforms the pointer-generator network [12] on the CNN/Daily Mail dataset.

#### D. Other Studies

So far, we primarily focused on the pointer-generator network, training neural networks with RL algorithms, and CNN based seq2seq architectures. There are many other studies that aim to improve the performance of seq2seq models for the task of abstractive text summarization from different perspectives and broaden their applications.

1) *Network Structure and Attention*: The first way to boost the performance of seq2seq models is to design better network structures. Zhou *et al.* [74] introduced an information filter, namely, a selective gate network between the encoder and decoder. This model can control the information flow from the encoder to the decoder via constructing a second level representation of the source texts with the gate network. Zeng *et al.* [38] introduced a read-again mechanism to improve the quality of the representations of the source texts. Tan *et al.* [75] built a graph ranking model upon a hierarchical encoder-decoder framework, which enables the model to capture the salient information of the source documents and generate accurate, fluent and non-redundant summaries. Xia *et al.* [76] proposed a deliberation network that passes the decoding process multiple times (deliberation process), to polish the sequences generated by the previous decoding process. Li *et al.* [77] incorporated a sequence of variational auto-encoders [78, 79] into the decoder to capture the latent structure of the generated summaries.

2) *Extraction + Abstraction*: Another way to improve the abstractive text summarization is to make use of the salient information from the extraction process. Hsu *et al.* [80] proposed a unified framework that takes advantage of both extractive and abstractive summarization using a novel attention mechanism, which is a combination of the sentence-level attention (based on the extractive summarization [81]) and the word-level attention (based on the pointer-generator network [12]), inspired by the intuition that words in less attended sentences should have lower attention scores. Chen *et al.* [59] introduced a multi-step procedure, namely compression-paraphrase, for abstractive summarization, which first extracts salient sentences from documents and then rewrites them in order to get final summaries. Li *et al.* [82] introduced a guiding generation model, where the keywords in source texts is first retrieved with an extractive model [83]. Then, a guide network is applied to encode them to obtain the key information representations that will guide the summary generation process.

3) *Long Documents*: Compared to short articles and texts with moderate lengths, there are many challenges that arise in long documents, such as difficulty in capturing the salient information [84]. Nallapati *et al.* [14] proposed a hierarchical attention model to capture hierarchical structures of long

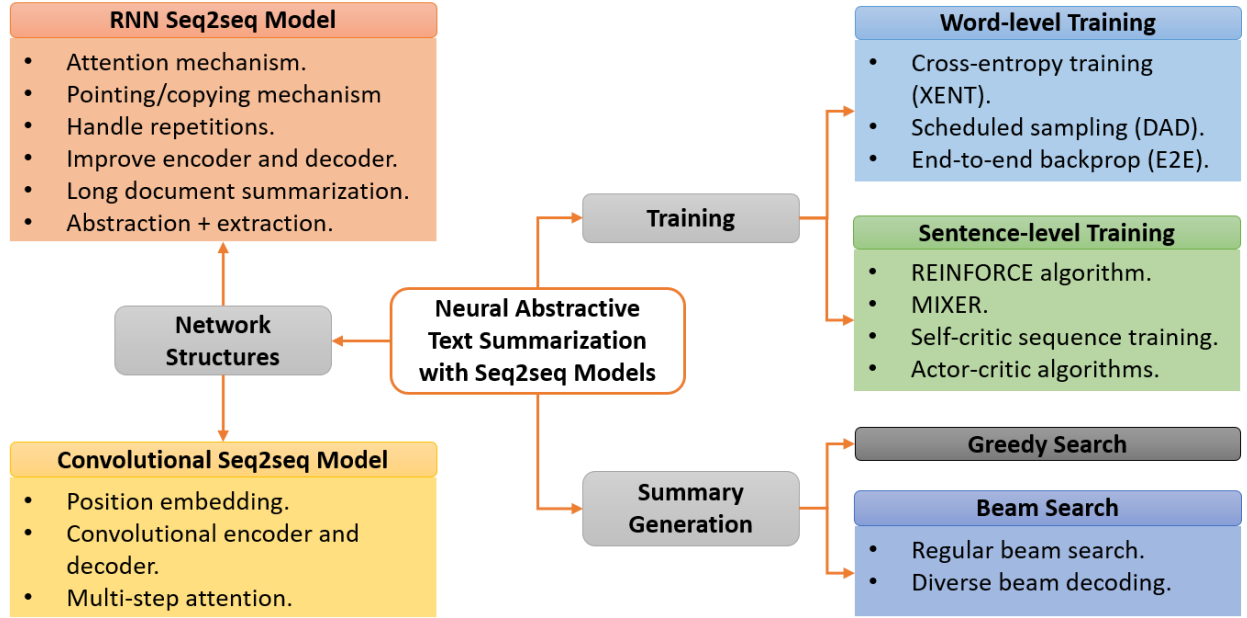


Fig. 1: An overall taxonomy of topics on seq2seq models for neural abstractive text summarization.

documents. To make models scale-up to very long sequences, Ling *et al.* [55] introduced a coarse-to-fine attention mechanism, which hierarchically reads and attends long documents<sup>2</sup>. By stochastically selecting chunks of texts during training, this approach can scale linearly with the number of chunks instead of the number of tokens. Cohan *et al.* [84] proposed a discourse-aware attention model which has a similar idea to that of a hierarchical attention model. Their model was applied to two large-scale datasets of scientific papers, i.e., arXiv and PubMed datasets. Tan *et al.* [75] introduced a graph-based attention model which is built upon a hierarchical encoder-decoder framework where the pagerank algorithm [85] was used to calculate saliency scores of sentences.

4) *Multi-Task Learning*: Multi-task learning has become a promising research direction for this problem since it allows seq2seq models to handle different tasks. Pasunuru *et al.* [86] introduced a multi-task learning framework, which incorporates knowledge from an entailment generation task into the abstractive text summarization task by sharing decoder parameters. They further proposed a novel framework [87] that is composed of two auxiliary tasks, i.e., question generation and entailment generation, to improve their model for capturing the saliency and entailment for the abstractive text summarization. In their model, different tasks share several encoder, decoder and attention layers. Mccann *et al.* [88] introduced a Natural Language Decathlon (decaNLP<sup>3</sup>), a challenge that spans ten different tasks, including question-answering, machine translation, summarization, and so on. They also proposed a multitask question answering network that can jointly learn all tasks without task-specific modules or parameters, since all tasks are mapped to the same framework of question-answering over a given context.

5) *Beam Search*: Beam search algorithms have been commonly used in the decoding of different language generation tasks [22, 12]. However, the generated candidate-sequences are usually lacking in diversity [89]. In other words, top- $K$  candidates are nearly identical, where  $K$  is size of a beam. Li *et al.* [90] replaced the log-likelihood objective function in the neural probabilistic language model [48] with Maximum Mutual Information (MMI) [91] in their neural conversation models to remedy the problem. This idea has also been applied to neural machine translation (NMT) [92] to model the bi-directional dependency of source and target texts. They further proposed a simple yet fast decoding algorithm that can generate diverse candidates and has shown performance improvement on the abstractive text summarization task [93]. Vijayakumar *et al.* [94] proposed generating diverse outputs by optimizing for a diversity-augmented objective function. Their method, referred to as Diverse Beam Search (DBS) algorithm, has been applied to image captioning, machine translation, and visual question-generation tasks. Cibils *et al.* [95] introduced a meta-algorithm that first uses DBS to generate summaries, and then, picks candidates according to maximal marginal relevance [96] under the assumption that the most useful candidates should be close to the source document and far away from each other. The proposed algorithm has boosted the performance of the pointer-generator network on CNN/Daily Mail dataset.

Despite many research papers that are published in the area of neural abstractive text summarization, there are few survey papers [97, 98, 99] that provide a comprehensive study. In this paper, we systematically review current advances of seq2seq models for the abstractive text summarization task from various perspectives, including network structures, training strategies, and sequence generation. In addition to a literature survey, we also implemented some of these methods

<sup>2</sup>A document is split into many chunks of texts.

<sup>3</sup><https://github.com/salesforce/decaNLP>

in an open-source library, namely NATS<sup>4</sup>. Extensive set of experiments have been conducted on various benchmark text summarization datasets in order to examine the importance of different network components. The main contributions of this paper can be summarized as follows:

- Provide a comprehensive literature survey of current advances of seq2seq models with an emphasis on the abstractive text summarization.
- Conduct a detailed review of the techniques used to tackle different challenges in RNN encoder-decoder architectures.
- Review different strategies for training seq2seq models and approaches for generating summaries.
- Provide an open-source library, which implements some of these models, and systematically investigate the effects of different network elements on the summarization performance.

The rest of this paper is organized as follows: An overall taxonomy of topics on seq2seq models for neural abstractive text summarization is shown in Fig. 1. A comprehensive list of papers published till date on the topic of neural abstractive text summarization have been summarized in Table I and II. In Section II, we introduce the basic seq2seq framework along with its extensions, including attention mechanism, pointing/copying mechanism, repetition handling, improving encoder or decoder, summarizing long documents and combining with extractive models. Section III summarizes different training strategies, including word-level training methods, such as cross-entropy training, and sentence-level training with RL algorithms. In Section IV, we discuss generating summaries using the beam search algorithm and various other diverse beam decoding algorithms. Section V briefly introduces the convolutional seq2seq model and its application to the abstractive text summarization. In Section VI, we present details of our implementations and discuss our experimental results on the CNN/Daily Mail, Newsroom [100], and Bytecup<sup>5</sup> datasets. We conclude this survey in Section VII.

## II. THE RNN ENCODER-DECODER FRAMEWORK

In this section, we review different encoder-decoder models for the neural abstractive text summarization. We will start with the basic RNN seq2seq framework and attention mechanism. Then, we will describe more advanced network structures that can handle different challenges in the text summarization, such as repetition and out-of-vocabulary (OOV) words. We will highlight various existing problems and proposed solutions.

### A. Seq2seq Framework Basics

A vanilla seq2seq framework for the abstractive summarization is composed of an encoder and a decoder. The encoder reads a source article, denoted by  $x = (x_1, x_2, \dots, x_J)$ , and transforms it to hidden states  $h^e = (h_1^e, h_2^e, \dots, h_J^e)$ ; while the decoder takes these hidden states as the context input and outputs a summary  $y = (y_1, y_2, \dots, y_T)$ . Here,  $x_i$  and  $y_j$

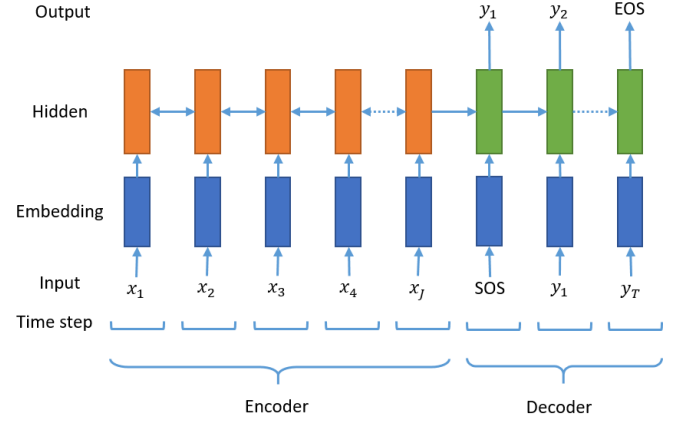


Fig. 2: The basic seq2seq model. SOS and EOS represent the start and end of a sequence, respectively.

are one-hot representations of the tokens in the source article and summary, respectively. We use  $J$  and  $T$  to represent the number of tokens (document length) of the original source document and the summary, respectively. A summarization task is defined as inferring a summary  $y$  from a given source article  $x$  using seq2seq models.

Encoders and decoders can be feed-forward networks, CNN [71, 67] or RNN. RNN architectures, especially long short term memory (LSTM) [32] and gated recurrent unit (GRU) [20], have been most widely adopted for seq2seq models. Fig. 2 shows a basic RNN seq2seq model with a bi-directional LSTM encoder and an LSTM decoder. The bi-directional LSTM is considered since it usually gives better document representations compared to a forward LSTM. The encoder reads a sequence of input tokens  $x$  and turns them into a sequences of hidden states  $h = (h_1, h_2, h_3, \dots, h_J)$  with following updating algorithm:

$$\begin{aligned}
 i_t &= \sigma(W_{ii}E_{x_{t-1}} + b_{ii} + W_{hi}h_{t-1} + b_{hi}) & (\text{Input Gate}) \\
 f_t &= \sigma(W_{if}E_{x_{t-1}} + b_{if} + W_{hf}h_{t-1} + b_{hf}) & (\text{Forget Gate}) \\
 o_t &= \sigma(W_{io}E_{x_{t-1}} + b_{io} + W_{ho}h_{t-1} + b_{ho}) & (\text{Output Gate}) \\
 g_t &= \tanh(W_{ig}E_{x_{t-1}} + b_{ig} + W_{hg}h_{t-1} + b_{hg}) \\
 c_t &= f_t c_{t-1} + i_t g_t \\
 h_t &= o_t \tanh(c_t)
 \end{aligned} \tag{1}$$

where weight matrices  $W$  and vector  $b$  are learnable parameters<sup>6</sup>,  $E_{x_t}$  denotes the word embeddings of token  $x_t$ , and  $c_t$  represents the cell states. Both  $h_t$  and  $c_t$  are initialized to 0. For the bi-directional LSTM, the input sequence is encoded as  $\overleftarrow{h}^e$  and  $\overrightarrow{h}^e$ , where the right and left arrows denote the forward and backward temporal dependencies, respectively. Superscript  $e$  is the shortcut notation used to indicate that it is for the encoder. During the decoding, the decoder takes the encoded representations of the source article (i.e., hidden and cell states  $\overleftarrow{h}_J^e, \overleftarrow{h}_1^e, \overleftarrow{c}_J^e, \overleftarrow{c}_1^e$ ) as the input and generates the summary  $y$ . In a simple encoder-decoder model, encoded vectors are used to

<sup>4</sup><https://github.com/tshi04/NATS>

<sup>5</sup><https://www.biendata.com/competition/bytecup2018/>

<sup>6</sup>In the rest of this paper, we will use  $W$  (weights) and  $b$  (bias) to represent the model parameters.

initialize hidden and cell states of the LSTM decoder. For example, we can initialize them as follows:

$$\begin{aligned} h_0^d &= \tanh(W_{e2d}(\vec{h}_J^e \oplus \overleftarrow{h}_1^e) + b_{e2d}) \\ c_0^d &= \vec{c}_J^e \oplus \overleftarrow{c}_1^e \end{aligned} \quad (2)$$

Here, superscript  $d$  denotes the decoder and  $\oplus$  is a concatenation operator. At each decoding step, we first update the hidden state  $h_t^d$  conditioned on the previous hidden states and input tokens, i.e.,

$$h_t^d = \text{LSTM}(h_{t-1}^d, E_{y_{t-1}}) \quad (3)$$

Hereafter, we will not explicitly express the cell states in the input and output of LSTM, since only hidden states are passed to other parts of the model. Then, the vocabulary distribution can be calculated as follows:

$$P_{\text{vocab},t} = \text{softmax}(W_{d2v}h_t^d + b_{d2v}), \quad (4)$$

where  $P_{\text{vocab},t}$  is a vector whose dimension is the size of the vocabulary  $\mathcal{V}$  and  $\text{softmax}(v_t) = \frac{\exp(v_t)}{\sum_r \exp(v_r)}$  for each element  $v_t$  of a vector  $v$ . Therefore, the probability of generating the target token  $w$  in the vocabulary  $\mathcal{V}$  is denoted as  $P_{\text{vocab},t}(w)$ .

This LSTM based encoder-decoder framework was the foundation of many neural abstractive text summarization models [14, 12, 17]. However, there are many problems with this model. For example, the encoder is not well trained via back propagation through time [101, 102], since the paths from encoder to the output are relatively far apart, which limits the propagation of gradient signals. The accuracy and human-readability of generated summaries is also very low with a lot of OOV words<sup>7</sup> and repetitions. The rest of this section will discuss different models that were proposed in the literature to resolve these issues for producing better quality summaries.

### B. Attention Mechanism

The attention mechanism has achieved great success and is commonly used in seq2seq models for different natural language processing (NLP) tasks [103], such as machine translation [23, 21], image captioning [104], and neural abstractive text summarization [14, 12, 17]. In an attention based encoder-decoder architecture (shown in Fig. 3), the decoder not only takes the encoded representations (i.e., final hidden and cell states) of the source article as input, but also selectively focuses on parts of the article at each decoding step. For example, suppose we want to compress the source input<sup>8</sup> “Kylian Mbappe scored two goals in four second-half minutes to send France into the World Cup quarter-finals with a thrilling 4-3 win over Argentina on Saturday.” to its short version “France beat Argentina 4-3 to enter quarter-finals.”. When generating the token “beat”, the decoder may need to attend “a thrilling 4-3 win” than other parts of the text. This attention can be achieved by an alignment mechanism [23],

<sup>7</sup>In the rest of this paper, we will use <unk>, i.e., unknown words, to denote OOV words.

<sup>8</sup><https://timesofindia.indiatimes.com/sports/football/fifa-world-cup/france-vs-argentina-live-score-fifa-world-cup-2018/articleshow/64807463.cms>

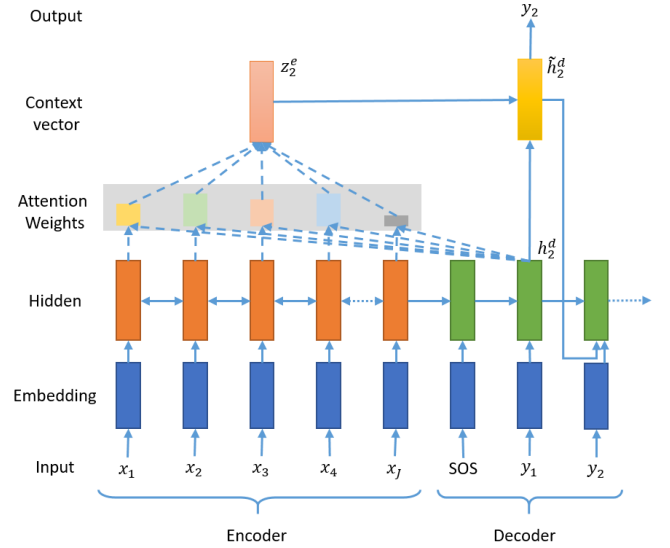


Fig. 3: An attention-based seq2seq model.

which first computes the attention distribution of the source tokens and then lets the decoder know where to attend to produce a target token. In the encoder-decoder framework depicted in Fig. 2 and 3, given all the hidden states of the encoder<sup>9</sup>, i.e.,  $h^e = (h_1^e, h_2^e, \dots, h_J^e)$  and the current decoder hidden state  $h_t^d$ , the attention distribution  $\alpha_t^e$  over the source tokens is calculated as follows:

$$\alpha_{tj}^e = \frac{\exp(s_{tj}^e)}{\sum_{k=1}^J \exp(s_{tk}^e)} \quad (5)$$

where the alignment score  $s_{tj}^e = s(h_j^e, h_t^d)$  is obtained by the content-based score function, which has three alternatives as suggested in [21]:

$$s(h_j^e, h_t^d) = \begin{cases} (h_j^e)^\top h_t^d & \text{dot} \\ (h_j^e)^\top W_{\text{align}} h_t^d & \text{general} \\ (v_{\text{align}})^\top \tanh(W_{\text{align}}(h_j^e \oplus h_t^d) + b_{\text{align}}) & \text{concat} \end{cases} \quad (6)$$

It should be noted that the number of additional parameters for ‘dot’, ‘general’ and ‘concat’ approaches are 0,  $|h_j^e| \times |h_t^d|$  and  $((|h_j^e| + |h_t^d|) \times |v_{\text{align}}| + 2 \times |v_{\text{align}}|)$ , respectively. Here  $|\cdot|$  represents the dimension of a vector. The ‘general’ and ‘concat’ are commonly used score functions in the abstractive text summarization [12, 17]. One of the drawbacks of ‘dot’ method is that it requires  $h_j^e$  and  $h_t^d$  to have the same dimension.

With the attention distribution, we can naturally define the source side context vector for the target word as

$$z_t^e = \sum_{j=1}^J \alpha_{tj}^e h_j^e \quad (7)$$

Together with the current decoder hidden state  $h_t^d$ , we get the attention hidden state [21]

$$\tilde{h}_t^d = W_z(z_t^e \oplus h_t^d) + b_z \quad (8)$$

<sup>9</sup> $h_j^e$  for the bi-directional LSTM is defined as the concatenation of  $\vec{h}_j^e$  and  $\overleftarrow{h}_j^e$ .



Finally, the vocabulary distribution is calculated by

$$P_{\text{vocab},t} = \text{softmax} \left( W_{d2v} \tilde{h}_t^d + b_{d2v} \right) \quad (9)$$

When  $t > 1$ , the decoder hidden state  $h_{t+1}^d$  is updated by

$$h_{t+1}^d = \text{LSTM} \left( h_t^d, E_{y_t} \oplus \tilde{h}_t^d \right) \quad (10)$$

where the input is concatenation of  $E_{y_t}$  and  $\tilde{h}_t^d$ .

### C. Pointing/Copying Mechanism

The pointing/copying mechanism [35] represents a class of approaches that generate target tokens by directly copying them from input sequences based on their attention weights. It can be naturally applied to the abstractive text summarization since summaries and articles can share the same vocabulary [12]. More importantly, it is capable to deal with out-of-vocabulary (OOV) words [14, 36, 37, 12]. A variety of studies have shown a boosting performance after incorporating the pointing/copying mechanism into the seq2seq framework [12, 17, 18]. In this section, we review several alternatives of this mechanism for the abstractive text summarization.

1) *Pointer Softmax* [37]: The basic architecture of pointer softmax is described as follows. It consists of three fundamental components: short-list softmax, location softmax and switching network. At decoding step  $t$ , a short-list softmax  $P_{\text{vocab},t}$  calculated by Eq. (9) is used to predict target tokens in the vocabulary. The location softmax gives locations of tokens that will be copied from the source article  $x$  to the target  $y_t$  based on attention weights  $\alpha_t^e$ . With these two components, a switching network is designed to determine whether to predict a token from the vocabulary or copy one from the source article if it is an OOV token. The switching network is a multilayer perceptron (MLP) with a sigmoid activation function, which estimates the probability  $p_{\text{gen},t}$  of generating tokens from the vocabulary based on the context vector  $z_t^e$  and hidden state  $h_t^d$  with

$$p_{\text{gen},t} = \sigma(W_{s,z} z_t^e + W_{s,h} h_t^d + b_s) \quad (11)$$

where  $p_{\text{gen},t}$  is a scalar and  $\sigma(a) = \frac{1}{1+\exp(-a)}$  is a sigmoid activation function. The final probability of producing the target token  $y_t$  is given by the concatenation of vectors  $p_{\text{gen},t} P_{\text{vocab},t}$  and  $(1 - p_{\text{gen},t}) \alpha_t^e$ .

2) *Switching Generator-Pointer* [14]: Similar to the switching network in pointer softmax [37], the switching generator-pointer is also equipped with a ‘switch’, which determines whether to generate a token from the vocabulary or point to one in the source article at each decoding step. The switch is explicitly modeled by

$$p_{\text{gen},t} = \sigma(W_{s,z} z_t^e + W_{s,h} h_t^d + W_{s,E} E_{y_{t-1}} + b_s) \quad (12)$$

If the switch is turned on, the decoder produces a word from the vocabulary with the distribution  $P_{\text{vocab},t}$  (see Eq. (9)). Otherwise, the decoder generates a pointer based on the attention distribution  $\alpha_t^e$  (see Eq. (5)), i.e.,  $p_j = \arg \max_{j \in \{1,2,\dots,J\}} \alpha_{tj}^e$ , where  $p_j$  is the position of the token in the source article. When a pointer is activated, embedding of the pointed token  $E_{x_j}$  will be used as an input for the next decoding step.

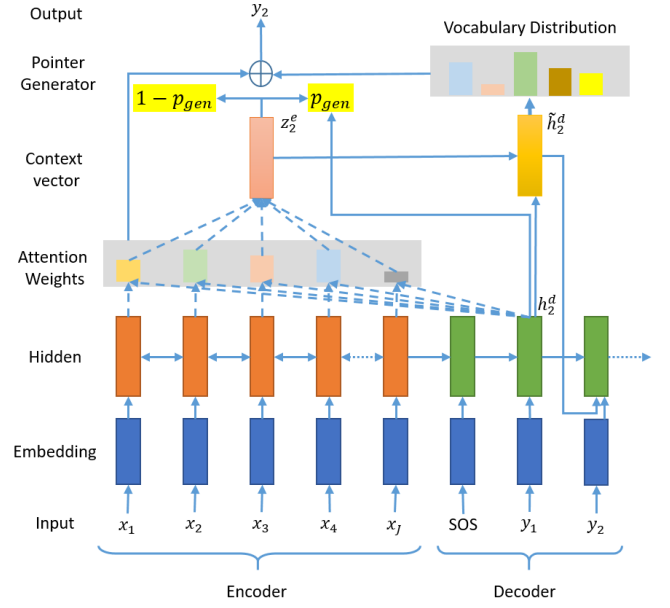


Fig. 4: The pointer-generator network.

3) *CopyNet* [36]: CopyNet has a differentiable network architecture and can be easily trained in an end-to-end manner. In this framework, the probability of generating a target token is a combination of the probabilities of two modes, i.e. generate-mode and copy-mode. First, CopyNet represents unique tokens in the vocabulary and source sequence by  $\mathcal{V}$  and  $\mathcal{X}$ , respectively, and builds an extended vocabulary  $\mathcal{V}_{\text{ext}} = \mathcal{V} \cup \mathcal{X} \cup \langle \text{unk} \rangle$ . Then, the vocabulary distribution over the extended vocabulary is calculated by

$$P_{\mathcal{V}_{\text{ext}}}(y_t) = P_g(y_t) + P_c(y_t) \quad (13)$$

where  $P_g$  and  $P_c$  are also defined on  $\mathcal{V}_{\text{ext}}$ , i.e.,

$$P_g(y_t) = \begin{cases} \frac{1}{Z} \exp \psi_g(y_t) & y_t \in \mathcal{V} \cup \langle \text{unk} \rangle \\ 0 & \text{otherwise} \end{cases} \quad (14)$$

$$P_c(y_t) = \begin{cases} \frac{1}{Z} \sum_{j:x_j=y_t} \exp \psi_c(x_j) & y_t \in \mathcal{X} \\ 0 & \text{otherwise} \end{cases} \quad (15)$$

Here,  $Z$  is a normalization factor shared by both the above equations.  $\psi_g(y_t)$  is calculated with

$$\psi_g(y_t) = W_{d2v} \tilde{h}_t^d + b_{d2v} \quad (16)$$

$\psi_c(x_j)$  is obtained by Eq. (6).

4) *Pointer-Generator Network* [12]: Pointer-generator network also has a differentiable network architecture (see Fig. 4). Similar to CopyNet [36], the vocabulary distribution over an extended vocabulary  $\mathcal{V}_{\text{ext}}$  is calculated by

$$P_{\mathcal{V}_{\text{ext}}}(y_t) = p_{\text{gen},t} P_g(y_t) + (1 - p_{\text{gen},t}) P_c(y_t) \quad (17)$$

where  $p_{\text{gen},t}$  is obtained by Eq. (12). Vocabulary distribution  $P_g(y_t)$  and attention distribution  $P_c(y_t)$  are defined as follows:

$$P_g(y_t) = \begin{cases} P_{\text{vocab},t}(y_t) & y_t \in \mathcal{V} \cup \langle \text{unk} \rangle \\ 0 & \text{otherwise} \end{cases} \quad (18)$$

and

$$P_c(y_t) = \begin{cases} \sum_{j:x_j=y_t} \alpha_{tj}^e & y_t \in \mathcal{X} \\ 0 & \text{otherwise} \end{cases} \quad (19)$$

The pointer-generator network has been used as the base model for many abstractive text summarization models (see Table I and II). Finally, it should be noted that  $p_{\text{gen},t} \in (0, 1)$  in CopyNet and pointer-generator network can be viewed as a “soft-switch” to choose between generation and copying, which is different from “hard-switch” (i.e.,  $p_{\text{gen},t} = 0, 1$ ) in pointer softmax and switching generator-pointer [37, 14, 17].

#### D. Repetition Handling

One of the critical challenges for attention based seq2seq models is that the generated sequences have repetitions, since the attention mechanism tends to ignore the past alignment information [105, 106]. For summarization and headline generation tasks, model generated summaries suffer from both word-level and sentence-level repetitions. The latter is specific to summaries which consist of several sentences [14, 12, 17], such as those in CNN/Daily Mail dataset [14] and Newsroom dataset [100]. In this section, we review several approaches that have been proposed to overcome the repetition problem.

1) *Temporal Attention* [14, 17]: Temporal attention method was originally proposed to deal with the attention deficiency problem in neural machine translation (NMT) [106]. Nallapati *et al.* [14] have found that it can also overcome the problem of repetition when generating multi-sentence summaries, since it prevents the model from attending the same parts of a source article by tracking the past attention weights. More formally, given the attention score  $s_{tj}^e$  in Eq. (6), we can first define a temporal attention score as [17]:

$$s_{tj}^{\text{temp}} = \begin{cases} \exp(s_{tj}^e) & \text{if } t = 1 \\ \frac{\exp(s_{tj}^e)}{\sum_{k=1}^{t-1} \exp(s_{kj}^e)} & \text{otherwise} \end{cases} \quad (20)$$

Then, attention distribution is calculated with

$$\alpha_{tj}^{\text{temp}} = \frac{s_{tj}^{\text{temp}}}{\sum_{k=1}^J s_{tk}^{\text{temp}}} \quad (21)$$

Given the attention distribution, the context vector (see Eq. (7)) is rewritten as

$$z_t^e = \sum_{j=1}^J \alpha_{tj}^{\text{temp}} h_j^e. \quad (22)$$

It can be seen from Eq. (20) that, at each decoding step, the input tokens which have been highly attended will have a lower attention score via the normalization in time dimension. As a result, the decoder will not repeatedly attend the same part of the source article.

2) *Intra-decoder Attention* [17]: Intra-decoder attention is another technique to handle the repetition problem for long-sequence generations. Compared to the regular attention based models, it allows a decoder to not only attend tokens in a source article but also keep track of the previously decoded tokens in a summary, so that the decoder will not repeatedly produce the same information.

For  $t > 1$ , intra-decoder attention scores, denoted by  $s_{t\tau}^d$ , can be calculated in the same manner as the attention score  $s_{tj}^e$ <sup>10</sup>. Then, the attention weight for each token is expressed as

$$\alpha_{t\tau}^d = \frac{\exp(s_{t\tau}^d)}{\sum_{k=1}^{t-1} \exp(s_{tk}^d)} \quad (23)$$

With attention distribution, we can calculate the decoder-side context vector by taking linear combination of the decoder hidden states, i.e.,  $h_{<t}^d$ , as

$$z_t^d = \sum_{\tau=1}^{t-1} \alpha_{t\tau}^d h_{\tau}^d \quad (24)$$

The decoder-side and encoder-side context vector will be both used to calculate the vocabulary distribution.

3) *Coverage* [12]: The coverage model was first proposed for the NMT task [105] to address the problems of the standard attention mechanism which tends to ignore the past alignment information. Recently, See *et al.* [12] introduced the coverage mechanism to the abstractive text summarization task. In their model, they first defined a coverage vector  $u_t^e$  as the sum of attention distributions of the previous decoding steps, i.e.,

$$u_t^e = \sum_j^{t-1} \alpha_{tj}^e \quad (25)$$

Thus, it contains the accumulated attention information on each token in the source article during the previous decoding steps. The coverage vector will then be used as an additional input to calculate the attention score

$$s_{tj}^e = (v_{\text{align}})^\top \tanh(W_{\text{align}}(h_j^d \oplus h_t^d \oplus u_t^e) + b_{\text{align}}) \quad (26)$$

As a result, the attention at current decoding time-step is aware of the attention during the previous decoding steps. Moreover, they defined a novel coverage loss to ensure that the decoder does not repeatedly attend the same locations when generating multi-sentence summaries. Here, the coverage loss is defined as

$$\text{covloss}_t = \sum_j \min(\alpha_{tj}^e, u_{tj}^e) \quad (27)$$

which is upper bounded by 1.

4) *Distraction* [107]: The coverage mechanism has also been used in [107] (known as *distraction*) for the document summarization task. In addition to the distraction mechanism over the attention, they also proposed a distraction mechanism over the encoder context vectors. Both mechanisms are used to prevent the model from attending certain regions of the source article repeatedly. Formally, given the context vector at current decoding step  $z_t^e$  and all historical context vectors  $(z_1^e, z_2^e, \dots, z_{t-1}^e)$  (see Eq. (7)), the distracted context vector  $z_t^{e,\text{dist}}$  is defined as

$$z_t^{e,\text{dist}} = \tanh(W_{\text{dist},z} z_t^e - W_{\text{hist},z} \sum_j^{t-1} z_j^e) \quad (28)$$

where both  $W_{\text{dist},z}$  and  $W_{\text{hist},z}$  are diagonal parameter matrices.

<sup>10</sup>We have to replace  $h_j^e$  with  $h_{\tau}^d$  in Eq. (6), where  $\tau \in \{1, \dots, t-1\}$ .



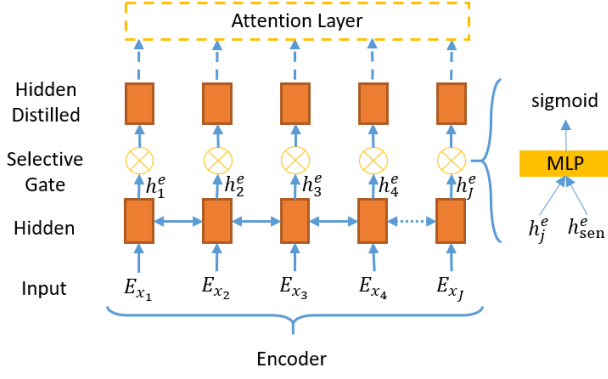


Fig. 5: An illustration of the selective encoder.

### E. Improving Encoded Representations

Although LSTM and bi-directional LSTM encoders<sup>11</sup> have been commonly used in the seq2seq models for the abstractive text summarization [14, 12, 17], representations of the source articles are still believed to be sub-optimal. In this section, we review some approaches that aim to improve the encoding process.

1) *Selective Encoding* [74]: The selective encoding model was proposed for the abstractive sentence summarization task [74]. Built upon an attention based encoder-decoder framework, it introduces a selective gate network into the encoder for the purpose of distilling salient information from source articles. A second layer representation, namely, distilled representation, of a source article is constructed over the representation of the first LSTM layer (a bi-directional GRU encoder in this work.). Formally, the distilled representation of each token in the source article is defined as

$$h_{\text{sel},j}^e = \text{gate}_{\text{sel},j} \times h_j^e \quad (29)$$

where  $\text{gate}_{\text{sel},j}$  denotes the selective gate for token  $x_j$  and is calculated as follows:

$$\text{gate}_{\text{sel},j} = \sigma(W_{\text{sel},h} h_j^e + W_{\text{sel},\text{sen}} h_{\text{sen}}^e + b_{\text{sel}}) \quad (30)$$

where  $h_{\text{sen}}^e = \vec{h_j^e} \oplus \overleftarrow{h_1^e}$ . The distilled representations are then used for the decoding. Such a gate network can control information flow from an encoder to a decoder and can also select salient information, therefore, it boosts the performance of the sentence summarization task [74].

2) *Read-Again Encoding* [38]: Intuitively, read-again mechanism is motivated by human readers who read an article several times before writing a summary. To simulate this cognitive process, a read-again encoder reads a source article twice and outputs two-level representations. In the first read, an LSTM encodes tokens and the article as  $(h_1^{e,1}, h_2^{e,1}, \dots, h_j^{e,1})$  and  $h_{\text{sen}}^{e,1} = h_j^{e,1}$ , respectively. In the second read, we use another LSTM to encode the source text based on the outputs of the first read. Formally, the encoder hidden state of the second read  $h_j^{e,2}$  is updated by

$$h_j^{e,2} = \text{LSTM}(h_{j-1}^{e,2}, E_{x_j} \oplus h_j^{e,1} \oplus h_{\text{sen}}^{e,1}) \quad (31)$$

<sup>11</sup>GRU and bi-directional GRU are also often seen in abstractive summarization papers.

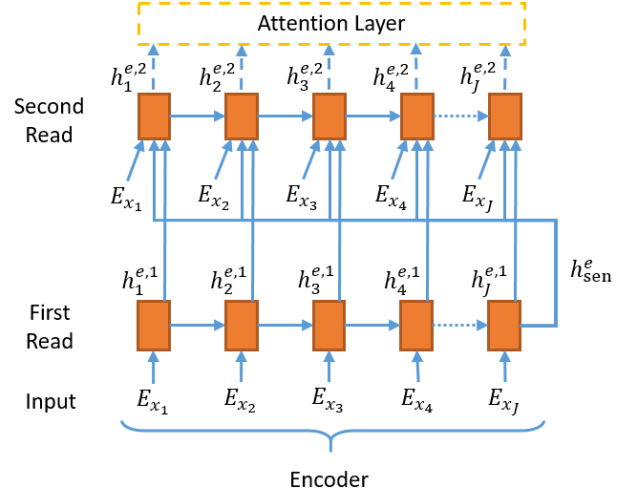


Fig. 6: An illustration of the read-again encoder.

The hidden states  $(h_1^{e,2}, h_2^{e,2}, \dots, h_j^{e,2})$  of the second read will be passed into decoders for summary generation.

### F. Improving Decoder

1) *Embedding Weight Sharing* [17]: Sharing the embedding weights with the decoder is a practical approach that can boost the performance since it allows us to reuse the semantic and syntactic information in an embedding matrix during summary generation [108, 17]. Suppose the embedding matrix is represented by  $W_{\text{emb}}$ , we can formulate the matrix used in the summary generation (see Eq. (9)) as follows:

$$W_{\text{d2v}} = \tanh(W_{\text{emb}}^\top \cdot W_{\text{proj}}) \quad (32)$$

By sharing model weights, the number of parameters is significantly less than a standard model since the number of parameters for  $W_{\text{proj}}$  is  $|h_j^e| \times |h_t^d|$ , while that for  $W_{\text{d2v}}$  is  $|h_t^d| \times |\mathcal{V}|$ , where  $|h|$  represents the dimension of vector  $h$  and  $|\mathcal{V}|$  denotes size of the vocabulary.

2) *Deliberation* [76]: When a human writes a summary, they usually first create a draft and then polish it based on the global context. Inspired by this polishing process, Xia *et al.* [76] proposed a deliberation network for sequence generation tasks. A deliberation network can have more than one decoder<sup>12</sup>. The first one is similar to the decoder presents in the basic seq2seq model described in Fig. 2. Let us denote the encoder hidden states by  $(h_1^e, h_2^e, \dots, h_j^e)$  and the first-pass decoder hidden states by  $(h_1^{d,1}, h_2^{d,1}, \dots, h_T^{d,1})$ . During the decoding, the second-pass decoder, which is used to polish the draft written in the first-pass, attends both encoder and the first-pass decoder. Therefore, we obtain two context vectors  $z_t^e = \sum_{j=1}^J \alpha_{tj}^e h_j^e$  and  $z_t^{d,1} = \sum_{\tau=1}^T \alpha_{t\tau}^{d,1} h_\tau^{d,1}$  at time step  $t$ , where  $\alpha_{tj}^e$  and  $\alpha_{t\tau}^{d,1}$  are attention weights. As we can see, two context vectors capture global information of the encoded article and the sequence generated by the first-pass decoder.

<sup>12</sup>There are two decoders in this paper

The second-pass decoder will take them as input and update the next hidden states with

$$h_{t+1}^{d,2} = \text{LSTM}(h_t^{d,2}, E_{y_t} \oplus z_t^e \oplus z_t^{d,1}) \quad (33)$$

Finally, the vocabulary distribution at decoding step  $t$  is calculated by

$$P_{\text{vocab},t} = \text{softmax}(W_{\text{d2v}}h_t^{d,2} + b_{\text{d2v}}) \quad (34)$$

The deliberation network has also boosted the performance of seq2seq models in NMT and abstractive text summarization tasks.

### 3) Deep Recurrent Generative Decoder (DRGD) [77]:

Conventional encoder-decoder models calculate hidden states and attention weights in an entirely deterministic fashion, which limits the capability of representations and results in low quality summaries. Incorporating variational auto-encoders (VAEs) [78, 79] into the encoder-decoder framework provides a practical solution for this problem. Inspired by the variational RNN proposed in [109] to model the highly structured sequential data, Li *et al.* [77] introduced a seq2seq model with DRGD that aims to capture latent structure information of summaries and improve the summarization quality. This model employs GRU as the basic recurrent model for both encoder and decoder. However, to be consistent with this survey paper, we will explain their ideas using LSTM instead.

There are two LSTM layers to calculate the decoder hidden state  $h_t^d$ . At the decoding step  $t$ , the first layer hidden state  $h_t^{d,1}$  is updated by  $h_t^{d,1} = \text{LSTM}^1(h_{t-1}^{d,1}, E_{y_{t-1}})$ . Then, the attention weights  $\alpha_{tj}^e$  and the context vector  $z_t^{d,1}$  are calculated with the encoder hidden state  $h^e$  and the first layer decoder hidden state  $h_t^{d,1}$  using Eqs. (5), (6) and (7). For the second layer, the hidden state  $h_t^{d,2}$  is updated with  $h_t^{d,2} = \text{LSTM}^2(h_{t-1}^{d,2}, E_{y_{t-1}} \oplus z_t^{d,1})$ . Finally, the decoder hidden state is obtained by  $h_t^d = h_t^{d,1} \oplus h_t^{d,2}$ , where  $h^d$  is also referred to as the deterministic hidden state.

VAE is incorporated into the decoder to capture latent structure information of summaries which is represented by a multivariate Gaussian distribution. By using a reparameterization trick [79, 110], latent variables can be first expressed as

$$\xi_t = \mu_t + \eta_t \otimes \epsilon, \quad (35)$$

where the noise variable  $\epsilon \sim \mathcal{N}(0, I)$ , and Gaussian parameters  $\mu_t$  and  $\eta_t$  in the network are calculated by

$$\mu_t = W_{\text{vae},\mu}h_t^{\text{enc}} + b_{\text{vae},\mu} \quad (36)$$

$$\log(\eta_t^2) = W_{\text{vae},\eta}h_t^{\text{enc}} + b_{\text{vae},\eta} \quad (37)$$

where  $h_t^{\text{enc}}$  is a hidden vector of the encoding process of the VAE and defined as

$$h_t^{\text{enc}} = \sigma(W_{\text{enc},\xi}\xi_{t-1} + W_{\text{enc},y}E_{y_{t-1}} + W_{\text{enc},h}h_{t-1}^d + b_{\text{enc}}) \quad (38)$$

With the latent structure variables  $\xi_t$ , the output hidden states  $h_t^{\text{dec}}$  can be formulated as

$$h_t^{\text{dec}} = \tanh(W_{\text{dec},\xi}\xi_t + W_{\text{dec},h}h_t^{d,2} + b_{\text{dec}}) \quad (39)$$

Finally, the vocabulary distribution is calculated by

$$P_{\text{vocab},t} = \text{softmax}(W_{\text{d2v}}h_t^{\text{dec}} + b_{\text{d2v}}) \quad (40)$$

We primarily focused on the network structure of DRGD in this section. The details of VAE and its derivations can be found in [78, 77, 79, 110]. In DRGD, VAE is incorporated into the decoder of a seq2seq model, more recent works have also used VAE in the attention layer [111] and for the sentence compression task [40].

### G. Summarizing Long Document

Compared with sentence summarization, the abstractive summarization for very long documents has been relatively less investigated. Recently, attention based seq2seq models with pointing/copying mechanism have shown their power in summarizing long documents with 400 and 800 tokens [12, 17]. However, performance improvement primarily attributes to copying and repetition/redundancy avoiding techniques [14, 12, 17]. For very long documents, we need to consider several important factors to generate high quality summaries [75], such as saliency, fluency, coherence and novelty. Usually, seq2seq models combined with the beam search decoding algorithm can generate fluent and human-readable sentences. In this section, we review models that aim to improve the performance of long document summarization from the perspective of saliency.

Seq2seq models for long document summarization usually consists of an encoder with a hierarchical architecture which is used to capture the hierarchical structure of the source documents. The top-level salient information includes the important sentences [14, 75], chunks of texts [55], sections [84], and paragraphs [18], while the lower-level salient information represents keywords. Hereafter, we will use the term ‘chunk’ to represent the top-level information. Fig. 7 shows neural network structure of a hierarchical encoder, which first uses a word-encoder to encode tokens in a chunk for the chunk representation, and then use a chunk encoder to encode the chunks in a document for the document representation. In this paper, we only consider the single-layer forward LSTM<sup>13</sup> for both word and chunk encoders.

Suppose, the hidden states of chunk  $i$  and word  $j$  in this chunk are represented by  $h_i^{\text{chk}}$  and  $h_{ij}^{\text{wd}}$ . At decoding step  $t$ , we can calculate word-level attention weight  $\alpha_{ij}^{\text{wd},t}$  for the current decoder hidden state  $h_t^d$  as follows:

$$\alpha_{ij}^{\text{wd},t} = \frac{\exp(s_{ij}^{\text{wd},t})}{\sum_{k,l} \exp(s_{kl}^{\text{wd},t})} \quad (41)$$

At the same time, we can also calculate chunk-level attention weight  $\alpha_i^{\text{chk},t}$  as follows:

$$\alpha_i^{\text{chk},t} = \frac{\exp(s_i^{\text{chk},t})}{\sum_k \exp(s_k^{\text{chk},t})} \quad (42)$$

where both alignment scores  $s_{ij}^{\text{wd},t} = s^{\text{wd}}(h_{ij}^{\text{wd}}, h_t^d)$  and  $s_i^{\text{chk},t} = s^{\text{chk}}(h_i^{\text{chk}}, h_t^d)$  can be calculated using Eq. (6). In this section, we will review four different models that are based on the hierarchical encoder for the task of long document text summarization.

<sup>13</sup>The deep communicating agents model [18], which requires multiple layers of bi-directional LSTM, falls out of the scope of this survey.

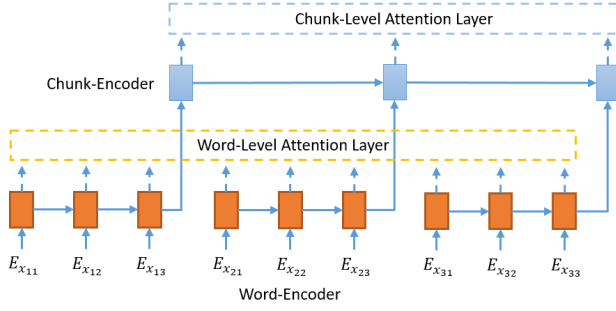


Fig. 7: A hierarchical encoder which first encodes tokens for the chunk representations and then encodes chunks for the document representation.

1) *Hierarchical Attention* [14]: The intuition behind a hierarchical attention is that words in less important chunks should be less attended. Therefore, with chunk-level attention distribution  $\alpha_i^{\text{chk},t}$  and word-level attention distribution  $\alpha_{ij}^{\text{wd},t}$ , we first calculate re-scaled word-level attention distribution by

$$\alpha_{ij}^{\text{scale},t} = \frac{\alpha_i^{\text{chk},t} \alpha_{ij}^{\text{wd},t}}{\sum_{k,l} \alpha_k^{\text{chk},t} \alpha_{kl}^{\text{wd},t}} \quad (43)$$

This re-scaled attention will then be used to calculate the context vector using Eq. (7), i.e.,

$$z_t^e = \sum_{i,j} \alpha_{ij}^{\text{scale},t} h_{ij}^{\text{wd}} \quad (44)$$

It should be noted that such hierarchical attention framework is different from the hierarchical attention network proposed in [112], where the chunk representation is obtained using

$$z_i^{\text{wd},t} = \sum_j \alpha_{ij}^{\text{wd},t} h_{ij}^{\text{wd}} \quad (45)$$

instead of the last hidden state of the word-encoder.

2) *Discourse-Aware Attention* [84]: The idea of the discourse-aware attention is similar to that of the hierarchical attention giving Eq. (43). The main difference between these two attention models is that the re-scaled attention distribution in the discourse-aware attention is calculated by

$$\alpha_{ij}^{\text{scale},t} = \frac{\exp(\alpha_i^{\text{chk},t} s_{ij}^{\text{wd},t})}{\sum_{k,l} \exp(\alpha_k^{\text{chk},t} s_{kl}^{\text{wd},t})} \quad (46)$$

3) *Coarse-to-Fine Attention* [55]: The coarse-to-fine (C2F) attention was proposed for computational efficiency. Similar to the hierarchical attention [14], the proposed model also has both chunk-level attention and word-level attention. However, instead of using word-level hidden states in all chunks for calculating the context vector, the C2F attention method first samples a chunk  $i$  from the chunk-level attention distribution, and then calculates the context vector using

$$z_t^e = \sum_j \alpha_{ij}^{\text{scale},t} h_{ij}^{\text{wd}} \quad (47)$$

At the test time, the stochastic sampling of the chunks will be replaced by a greedy search.

4) *Graph-based Attention* [75]: The aforementioned hierarchical attention mechanism implicitly captures the chunk-level salient information, where the importance of a chunk is determined solely by its attention weight. In contrast, the graph-based attention framework allows us to calculate the saliency scores explicitly using the pagerank algorithm [85, 113] on a graph whose vertices and edges are chunks of texts and their similarities, respectively. Formally, at the decoding time-step  $t$ , saliency scores for all input chunks are obtained by

$$f^t = (1 - \lambda)(I - \lambda W^{\text{adj}}(t) D_{\text{adj}}^{-1}(t))^{-1} \chi_{\mathcal{T}} \quad (48)$$

where adjacent matrix  $W^{\text{adj}}$  (similarity of chunks) is calculated by

$$W_{ij}^{\text{adj}} = h_i^{\text{chk}} W_{\text{par}}^{\text{chk}} h_j^{\text{chk}} \quad (49)$$

$D_{\text{adj}}$  is a diagonal matrix with its  $(i, i)$ -element equal to the sum of the  $i^{\text{th}}$  column of  $W^{\text{adj}}$ .  $\lambda$  is a damping factor. The vector  $\chi_{\mathcal{T}}$  is defined as

$$\chi_{\mathcal{T},i} = \begin{cases} \frac{1}{|\mathcal{T}|} & i \in \mathcal{T} \\ 0 & \text{otherwise} \end{cases} \quad (50)$$

where  $\mathcal{T}$  is a topic (see [113, 75] for more details). Finally, the graph-based attention distribution over a chunk can be obtained by

$$\alpha_i^{\text{chk},t} = \frac{\max(f_i^t - f_i^{t-1}, 0)}{\sum_k (\max(f_k^t - f_k^{t-1}, 0))} \quad (51)$$

where  $f_i^0$  is initialized with 0. It can be seen that the graph-based attention mechanism will focus on chunks that rank higher than the previous decoding step, i.e.,  $f_i^t > f_i^{t-1}$ . Therefore, it provides an efficient way to select salient information from source documents.

## H. Extraction + Abstraction

Extractive summarization approaches usually show a better performance comparing to the abstractive approaches [12, 14, 13] especially with respect to ROUGE measures. One of the advantages of the extractive approaches is that they can summarize source articles by extracting salient snippets and sentences directly from these documents [81], while abstractive approaches rely on word-level attention mechanism to determine the most relevant words to the target words at each decoding step. In this section, we review several studies that have attempted to improve the performance of the abstractive summarization by combining them with extractive models.

1) *Extractor + Pointer-Generator Network* [80]: This model proposes a unified framework that tries to leverage the sentence-level salient information from an extractive model and incorporate them into an abstractive model (a pointer-generator network). More formally, inspired by the hierarchical attention mechanism [14], they replaced the attention distribution  $\alpha_t^e$  in the abstractive model with a scaled version  $\alpha_t^{\text{scale}}$ , where the attention weights are expressed as follows:

$$\alpha_{tj}^{\text{scale}} = \frac{\alpha_{tj}^{\text{extra}} \alpha_{tj}^{\text{wd}}}{\sum_k \alpha_{tk}^{\text{extra}} \alpha_{tk}^{\text{wd}}} \quad (52)$$

TABLE I: An overview of different seq2seq models for the neural abstractive text summarization (2015-2017).

Year	Reference	Highlights	Framework	Training	Optimizer	Datasets				Metrics
						DUC	Gigaword	CNN/DM	Others	
2015	Rush <i>et al.</i> [15]	Attention Based Summarization (ABS)	Bag-of-words, Convolution, Attention → Neural Network Language Model (NNLM)	XENT	SGD	✓	✓	-	-	ROUGE
	lopyrev <i>et al.</i> [114]	Simple Attention	LSTM → LSTM	XENT	RMSProp	-	✓	-	-	BLEU
	Ranzato <i>et al.</i> [46]	Sequence-level Training	Elman, LSTM → Elman, LSTM	XENT, DAD, E2E, MIXER	SGD	-	✓	-	-	ROUGE, BLEU
2016	Chopra <i>et al.</i> [26]	Recurrent Attentive Summarizer	Convolution Encoder, Attentive Encoder → Elman, LSTM	XENT	SGD	✓	✓	-	-	ROUGE
	Nallapati <i>et al.</i> [14]	Switch Generator-Pointer, Temporal-Attention, Hierarchical-Attention	RNN, Feature-rich Encoder → RNN	XENT	Adadelata	✓	✓	✓	-	ROUGE
	Miao <i>et al.</i> [40]	Auto-encoding Sentence Compression, Forced-Attention Sentence Compression, Pointer Network	Encoder → Compressor → Decoder	XENT+RL	Adam	-	✓	-	-	ROUGE
	Chen <i>et al.</i> [107]	Distraction	GRU → GRU	XENT	Adadelata	-	-	-	CNN, LCSTS	ROUGE
	Gulcehre <i>et al.</i> [37]	Pointer softmax	GRU → GRU	XENT	Adadelata	-	✓	-	-	ROUGE
	Gu <i>et al.</i> [36]	CopyNet	GRU → GRU	XENT	SGD	-	-	-	LCSTS	ROUGE
	Zeng <i>et al.</i> [38]	Read-again, Copy Mechanism	LSTM/GRU/Hierarchical read-again encoder → LSTM	XENT	SGD	✓	✓	-	-	ROUGE
	Li <i>et al.</i> [93]	Diverse Beam Decoding	LSTM → LSTM	RL	SGD	-	-	-	-	ROUGE
	Takase <i>et al.</i> [115]	Abstract Meaning Representation (AMR) based on ABS.	Attention-based AMR encoder → NNLM	XENT	SGD	✓	✓	-	-	ROUGE
2017	See <i>et al.</i> [12]	Pointer-Generator Network, Coverage	LSTM → LSTM	XENT	Adadelata	-	-	✓	-	ROUGE, METER
	Paulus <i>et al.</i> [17]	A Deep Reinforced Model, Intra-temporal and Intra-decoder Attention, Weight Sharing	LSTM → LSTM	XENT + RL	Adam	-	-	✓	NYT	ROUGE, Human
	Zhou <i>et al.</i> [74]	Selective Encoding, Abstractive Sentence Summarization	GRU → GRU	XENT	SGD	✓	✓	-	MSR-ATC	ROUGE
	Xia <i>et al.</i> [76]	Deliberation Networks	LSTM → LSTM	XENT	Adadelata	-	✓	-	-	ROUGE
	Nema <i>et al.</i> [116]	Query-based, Diversity based Attention	GRU query encoder, document encoder → GRU	XENT	Adam	-	-	-	Debate-pedia	ROUGE
	Tan <i>et al.</i> [75]	Graph-based Attention	Hierarchical Encoder → LSTM	XENT	Adam	-	-	✓	CNN, DailyMail	ROUGE
	Ling <i>et al.</i> [55]	Coarse-to-fine Attention	LSTM → LSTM	RL	SGD	-	-	✓	-	ROUGE, PPL
	Zhang <i>et al.</i> [56]	Sentence Simplification, Reinforcement Learning	LSTM → LSTM	RL	Adam	-	-	-	Newsela, WikiSmall, WikiLarge	BLEU, FKGL, SARI
	Li <i>et al.</i> [77]	Deep Recurrent Generative Decoder (DRGD)	GRU → GRU, VAE	XENT, VAE	Adadelata	✓	✓	-	LCSTS	ROUGE
	Liu <i>et al.</i> [117]	Adversarial Training	Pointer-Generator Network	GAN	Adadelata	-	-	✓	-	ROUGE, Human
	Pasunuru <i>et al.</i> [86]	Multi-Task with Entailment Generation	LSTM document encoder and premise Encoder → LSTM Summary and Entailment Decoder	Hybrid-Objective	Adam	✓	✓	-	SNLI	ROUGE, METEOR, BLEU, CIDEr-D
	Gehring <i>et al.</i> [67]	Convolutional Seq2seq, Position Embeddings, Gated Linear Unit, Multi-step Attention	CNN → CNN	XENT	Adam	✓	✓	-	-	ROUGE
	Fan <i>et al.</i> [72]	Convolutional Seq2seq, Controllable	CNN → CNN	XENT	Adam	✓	-	✓	-	ROUGE, Human

TABLE II: An overview of different seq2seq models for the neural abstractive text summarization (2018).

Year	Reference	Highlights	Framework	Training	Optimizer	Datasets				Metrics
						DUC	Gigaword	CNN/DM	Others	
2018	Celikyilmaz <i>et al.</i> [18]	Deep Communicating Agents, Semantic Cohesion Loss	LSTM $\rightarrow$ LSTM	Hybrid-Objective	Adam	-	-	✓	NYT	ROUGE, Human
	Chen <i>et al.</i> [59]	Reinforce-Selected Sentence Rewriting	LSTM Encoder $\rightarrow$ Extractor $\rightarrow$ Abstractor	XENT + RL	SGD	✓	-	✓	-	ROUGE, Human
	Hsu <i>et al.</i> [80]	Abstraction + Extraction, Inconsistency Loss	Extractor: GRU. Abstractor: Pointer-generator Network	Hybrid-Objective + RL	Adadelata	-	-	✓	-	ROUGE, Human
	Li <i>et al.</i> [58]	Actor-Critic	GRU $\rightarrow$ GRU	RL	Adadelata	✓	✓	-	LCSTS	ROUGE
	Li <i>et al.</i> [82]	Abstraction + Extraction, Key Information Guide Network (KIGN)	KIGN: LSTM. Framework: Pointer-Generator Network	XENT	Adadelata	-	-	✓	-	ROUGE
	Lin <i>et al.</i> [118]	Global Encoding, Convolutional Gated Unit	LSTM $\rightarrow$ LSTM	XENT	Adam	-	✓	-	LCSTS	ROUGE
	Pasunuru <i>et al.</i> [57]	Multi-Reward Optimization for RL: ROUGE, Saliency and Entailment.	LSTM $\rightarrow$ LSTM	RL	Adam	✓	-	✓	SNLI, MultiNLI, SQuAD	ROUGE, Human
	Song <i>et al.</i> [41]	Structured-Infused Copy Mechanisms	Pointer-Generator Network	Hybrid-Objective	Adam	-	✓	-	-	ROUGE, Human
	Cohan <i>et al.</i> [84]	Discourse Aware Attention	Hierarchical RNN LSTM Encoder $\rightarrow$ LSTM	XENT	Adagrad	-	-	-	PubMed, arXiv	ROUGE
	Guo <i>et al.</i> [87]	Multi-Task Summarization with Entailment and Question Generation	Multi-Task Encoder-Decoder Framework	Hybrid-Objective	Adam	✓	✓	✓	SQuAD, SNLI	ROUGE, METEOR
	Cibils <i>et al.</i> [95]	Diverse Beam Search, Plagiarism and Extraction Scores	Pointer-Generator Network	XENT	Adagrad	-	-	✓	-	ROUGE
	Wang <i>et al.</i> [119]	Topic Aware Attention	CNN $\rightarrow$ CNN	RL	-	✓	✓	-	LCSTS	ROUGE
	Kryściński <i>et al.</i> [120]	Improve Abstraction	LSTM Encoder $\rightarrow$ Decoder: Contextual Model and Language Model	XENT + RL	Asynchronous Gradient Descent Optimizer	-	-	✓	-	ROUGE, Novel n-gram Test, Human
	Gehrmann <i>et al.</i> [43]	Bottom-up Attention, Abstraction + Extraction	Pointer-Generator Network	Hybrid-Objective	Adagrad	-	-	✓	NYT	ROUGE, %Novel
	Zhang <i>et al.</i> [121]	Learning to Summarize Radiology Findings	Pointer-Generator Network + Background Encoder	XENT	Adam	-	-	-	Radiology Reports	ROUGE
	Jiang <i>et al.</i> [42]	Closed-book Training	Pointer-Generator Network + Closed-book Decoder	Hybrid-Objective + RL	Adam	✓	-	✓	-	ROUGE, METEOR
	Chung <i>et al.</i> [122]	Main Pointer Generator	Pointer-Generator Network + Document Encoder	XENT	Adadelata	-	-	✓	-	ROUGE
	Chen <i>et al.</i> [123]	Iterative Text Summarization	GRU encoder, GRU decoder, iterative unit	Hybrid-Objective	Adam	✓	-	✓	-	ROUGE

Here,  $\alpha_{tj}^{\text{extra}}$  is the sentence-level salient score of the sentence at word position  $j$  and decoding step  $t$ . Different from [14], the salient scores (sentence-level attention weights) are obtained from another deep neural network known as extractor [80].

During training, in addition to cross-entropy and coverage loss used in the pointer-generator network, this paper also proposed two other losses, i.e., *extractor loss* and *inconsistency loss*. The *extractor loss* is used to train the extractor and is defined as follows:

$$L_{\text{ext}} = -\frac{1}{N} \sum_{n=1}^N g_n \log \beta_n + (1 - g_n) \log(1 - \beta_n) \quad (53)$$

where  $g_n$  is the ground truth label for the  $n^{\text{th}}$  sentence and  $N$  is the total number of sentences. The *inconsistency loss* is

expressed as

$$L_{\text{inc}} = -\frac{1}{T} \sum_{t=1}^T \log\left(\frac{1}{|\mathcal{K}|} \sum_{j \in \mathcal{K}} \alpha_{tj}^e \alpha_{tj}^{\text{extra}}\right) \quad (54)$$

where  $\mathcal{K}$  is the set of the top- $k$  attended words and  $T$  is the total number of words in a summary. Intuitively, the inconsistency loss is used to ensure that the sentence-level attentions in the extractive model and word-level attentions in the abstractive model are consistent with each other. In other words, when word-level attention weights are high, the corresponding sentence-level attention weights should also be high.

2) *Key-Information Guide Network (KIGN)* [82]: This approach uses a guiding generation mechanism that leverages the key (salient) information, i.e., keywords, to guide decoding process. This is a two-step procedure. First, keywords are ex-

tracted from source articles using the TextRank algorithm [83]. Second, a KIGN encodes the key information and incorporates them into the decoder to guide the generation of summaries. Technically speaking, we can use a bi-directional LSTM to encode the key information and the output vector is the concatenation of hidden states, i.e.,  $h^{\text{key}} = h_N^{\text{key}} \oplus h_1^{\text{key}}$ , where  $N$  is the length of the key information sequence. Then, the alignment mechanism is modified as

$$s_{tj}^e = (v_{\text{align}})^{\top} \tanh(W_{\text{align}}^e h_j^e + W_{\text{align}}^d h_t^d + W_{\text{align}}^{\text{key}} h^{\text{key}}) \quad (55)$$

Similarly, the soft-switch in the pointer-generator network is calculated using

$$p_{\text{gen},t} = \sigma(W_{s,z} z_t^e + W_{s,h} h_t^d + W_{s,\text{key}} h^{\text{key}} + b_s) \quad (56)$$

3) *Reinforce-Selected Sentence Rewriting* [59]: Most models introduced in this survey are built upon the encoder-decoder framework [14, 12, 17], in which the encoder reads source articles and turns them into vector representations, and the decoder takes the encoded vectors as input and generates summaries. Unlike these models, the reinforce-selected sentence rewriting model [59] consists of two seq2seq models. The first one is an extractive model (*extractor*) which is designed to extract salient sentences from a source article, while the second is an abstractive model (*abstractor*) which paraphrases and compresses the extracted sentences into a short summary. The abstractor network is a standard attention-based seq2seq model with the copying mechanism for handling OOV words. For the extractor network, an encoder first uses a CNN to encode tokens and obtains representations of sentences, and then it uses an LSTM to encode the sentences and represent a source document. With the sentence-level representations, the decoder (another LSTM) is designed to recurrently extract salient sentences from the document using the pointing mechanism [35]. This model has achieved the state-of-the-art performance on CNN/Daily Mail dataset and was demonstrated to be computationally more efficient than the pointer-generator network [12].

### III. TRAINING STRATEGIES

In this section, we review different strategies to train the seq2seq models for abstractive text summarization. As discussed in [46], there are two categories of training methodologies, i.e., word-level and sequence-level training. The commonly used teacher forcing algorithm [44, 49] and cross-entropy training [48, 124] belong to the first category, while different RL-based algorithms [46, 53, 52] fall into the second. We now discuss the basic ideas of different training algorithms and their applications to seq2seq models for the text summarization. A comprehensive survey of deep RL for seq2seq models can be found in [47].

#### A. Word-Level Training

The word-level training for language models represents methodologies that try to optimize predictions of the next token [46]. For example, in the abstractive text summarization, given a source article  $x$ , a seq2seq model generates a summary

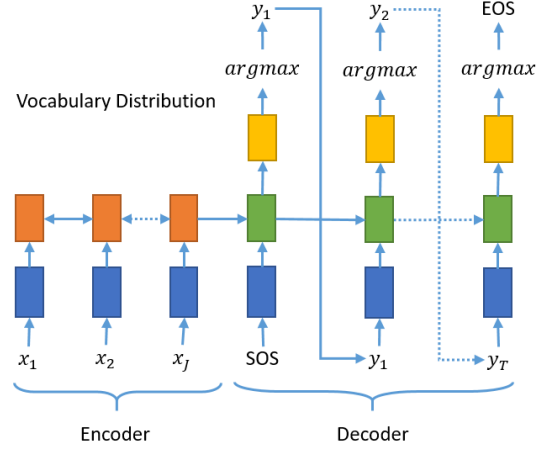


Fig. 8: Generation process with a greedy search.

$y$  with the probability  $P_{\theta}(y|x)$ , where  $\theta$  represents model parameters (e.g., weights  $W$  and bias  $b$ ). In a neural language model [48], this probability can be expanded to

$$P_{\theta}(y|x) = \prod_{t=1}^T P_{\theta}(y_t|y_{<t}, x), \quad (57)$$

where each multiplier  $P_{\theta}(y_t|y_{<t}, x)$ , known as likelihood, is a conditional probability of the next token  $y_t$  given all previous ones denoted by  $y_{<t} = (y_1, y_2, \dots, y_{t-1})$ . Intuitively, the text generation process can be described as follows. Starting with a special token ‘SOS’ (start of sequence), the model generates a token  $y_t$  at a time  $t$  with the probability  $P_{\theta}(y_t|y_{<t}, x) = P_{\text{vocab},t}(y_t)$ . This token can be obtained by a sampling method or a greedy search, i.e.,  $y_t = \arg \max_{y_t} P_{\text{vocab},t}$  (see Fig. 8). The generated token will then be fed into the next decoding step. The generation is stopped when the model outputs ‘EOS’ (end of sequence) token or when the length reaches a user defined maximum threshold. In this section, we review different approaches for learning model parameters, i.e.,  $\theta$ . We will start with the commonly used end-to-end training approach, i.e., cross-entropy training, and then move on to two different methods for avoiding the problem of exposure bias.

1) *Cross-Entropy Training (XENT)* [46]: To learn model parameters  $\theta$ , XENT maximizes the log-likelihood of observed sequences (ground-truth)  $\hat{y}_t = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_T)$ , i.e.,

$$\log P_{\theta}(\hat{y}|x) = \sum_{t=1}^T \log P_{\theta}(\hat{y}_t|\hat{y}_{<t}, x) \quad (58)$$

which is equivalent to minimizing the cross entropy (XE) loss,

$$\text{loss}_{\text{XE}} = -\log P_{\theta}(\hat{y}|x) \quad (59)$$

We show this training strategy in Fig. 9. The algorithm is also known as the teacher forcing algorithm [44, 49]. During training, it uses observed tokens (ground-truth) as input and aims to improve the probability of the next observed token at each decoding step. However, during testing, it relies on predicted tokens from the previous decoding step. This is



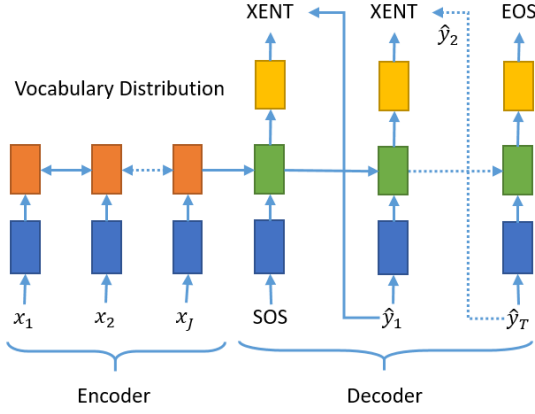


Fig. 9: Training with the teacher forcing algorithm.

the major difference between training and testing (see Fig. 8 and Fig. 9). Since the predicted tokens may not be the observed ones, this discrepancy will be accumulated over time and thus yields summaries that are very different from ground-truth summaries. This problem is known as exposure bias [46, 45, 44].

2) *Scheduled Sampling* [46, 45, 44]: Scheduled sampling algorithm, also known as Data As Demonstrator (DAD) [46, 45], has been proposed to solve the exposure bias problem. As shown in Fig. 10, during training, the input at each decoding step comes from a sampler which can decide whether it is a model generated token  $y_t$  from the last step or an observed token  $\hat{y}_t$  from training data. The sampling is based on a Bernoulli distribution

$$P_{\text{dad}}(y) = p_{\text{dad}}^{I(y=\hat{y}_t)} \cdot (1 - p_{\text{dad}})^{I(y=y_t)} \quad (60)$$

where  $p_{\text{dad}}$  is the probability of using a token from training data and  $I(y = y_t)$  is a binary indicator function. In the scheduled sampling algorithm,  $p_{\text{dad}}$  is an annealing/scheduling function and decreases with training time from 1 to 0. As suggested by Bengio *et al.* [44], scheduling function can take different forms, e.g.,

$$p_{\text{dad}} = \begin{cases} 1 - \alpha k & \text{linear decay} \\ \alpha^k & \text{exponential decay} \\ \frac{\alpha}{\alpha + \exp(k/\alpha)} & \text{inverse sigmoid decay} \end{cases} \quad (61)$$

where  $k$  is training step and  $\alpha$  is a parameter that guarantees  $p_{\text{dad}} \in [0, 1]$ . This strategy is often referred to as a curriculum learning algorithm [44, 56, 46].

The main intuition behind this algorithm is that, at the beginning stage, a model with random parameters cannot generate relevant/correct tokens, therefore, a decoder takes ground-truth tokens from training data as input. As the training proceeds, the model gradually reduces the probability of using ground-truth tokens. By the end of the training, the model assumes that it has been well trained and can generate reasonable tokens, thus, the decoder can completely rely on its own predictions [44].

3) *End-To-End Backprop (E2E)* [46]: This algorithm is another method that exposes a model to its own predictions

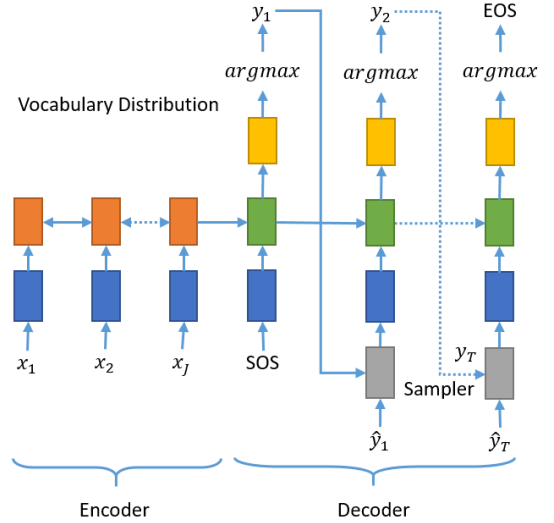


Fig. 10: Illustration of the scheduled sampling.

during training. At each decoding step, it still uses XENT to train the model parameters. However, the input is neither a ground truth token nor a model generated token. Instead, it is a fusion of top- $k$  tokens from the last decoding step, where  $k$  is a hyper-parameter. More specifically, the model first samples the top- $k$  tokens, denoted as  $(y_t^{\text{samp}_1}, y_t^{\text{samp}_2}, \dots, y_t^{\text{samp}_k})$ , from the vocabulary distribution  $P_{\text{vocab},t}$ . Then, it can re-scale their probabilities as follows:

$$P_{\text{samp},t}(y_t^{\text{samp}_i}) = \frac{P_{\text{vocab},t}(y_t^{\text{samp}_i})}{\sum_j P_{\text{vocab},t}(y_t^{\text{samp}_j})} \quad (62)$$

and obtain a vector in the embedding space by

$$E_{\text{samp},t} = \sum_{i=1}^k P_{\text{samp},t}(y_t^{\text{samp}_i}) E_{y_t^{\text{samp}_i}}. \quad (63)$$

This fused vector will be served as the input for the next decoding step. It should be noted that the proposed E2E algorithm also makes use of the DAD algorithm in practice [46], where a sampler is used to determine whether to take the fused vectors or the embeddings of the ground-truth tokens as the input.

## B. Sequence-Level Training

The sequence-level training with deep RL algorithms has recently received a lot of popularity in the area of neural abstractive text summarization [46, 56, 17, 47, 57, 58] due to its ability to incorporate any user-defined metrics, including non-differentiable ROUGE scores, to train neural networks. In this section, we review several different policy gradient algorithms that have been used to train abstractive text summarization models. For actor-critic algorithms and related work, the readers are encouraged to go through these publications [47, 58, 59].

In RL setting, generating a sequence of tokens in a summary can be considered as a sequential decision making process, where an encoder-decoder model is viewed as an agent, which

first reads a source article  $x$  and initializes its internal state (hidden and cell states for LSTM). At the decoding step  $t$ , it updates the state and takes an action  $y_t^\pi \in \mathcal{V}$  (i.e., picking a token from the vocabulary) according to a policy  $\pi = P_\theta(y_t|y_{<t}, x)$ . Here, the vocabulary is viewed as an action space. By the end of the decoding, it will produce a sequence of actions  $(y_1^\pi, y_2^\pi, \dots, y_T^\pi)$  and observe a reward  $R(y_1^\pi, y_2^\pi, \dots, y_T^\pi)$ , which is usually ROUGE scores [16] in the context of text summarization. Then, RL algorithms will be used to update the agent by comparing the action sequence based on current policy with the optimal action sequence (i.e., the ground-truth summary). In this section, we will start with the commonly used REINFORCE algorithm for training seq2seq models. Then, we will introduce the MIXER algorithm (which can improve the convergence rate and stability of the training) and the self-critic sequence training approach (which shows low variance for the gradient estimator).

1) *REINFORCE algorithm* [51, 46, 56, 55]: The goal of REINFORCE algorithm is to find parameters that maximize the expected rewards. Therefore, the loss function is defined as negative expected reward, i.e.,

$$L_\theta = -\mathbb{E}_{(y_{<T+1}^\pi) \sim \pi_\theta} R(y_{<T+1}^\pi) \quad (64)$$

where  $y_{<T+1}^\pi$  represents  $y_1^\pi, y_2^\pi, \dots, y_T^\pi$ . The above equation can also be rewritten as

$$L_\theta = - \sum_{y_{<T+1}^\pi \in \mathcal{Y}} \pi_\theta(y_{<T+1}^\pi) R(y_{<T+1}^\pi) \quad (65)$$

where  $\mathcal{Y}$  represents a set that contains all possible sequences. To optimize policy with respect to model parameters, we take the derivative of the loss function and obtain

$$\begin{aligned} \nabla_\theta L_\theta &= - \sum_{y_{<T+1}^\pi \in \mathcal{Y}} \nabla_\theta \pi_\theta(y_{<T+1}^\pi) R(y_{<T+1}^\pi) \\ &= - \sum_{y_{<T+1}^\pi \in \mathcal{Y}} \pi_\theta(y_{<T+1}^\pi) \nabla_\theta \log \pi_\theta(y_{<T+1}^\pi) R(y_{<T+1}^\pi) \end{aligned} \quad (66)$$

In the abstractive text summarization task, the policy is expressed as  $\pi_\theta = P_\theta(y_{<T+1}^\pi|x)$  and according to Eq. (57), the above equation can be expressed as follows:

$$\begin{aligned} \nabla_\theta L_\theta &= - \sum_{y_{<T+1}^\pi \in \mathcal{Y}} P_\theta(y_{<T+1}^\pi|x) \cdot \left[ \sum_{t=1}^T \nabla_\theta \log P_\theta(y_t^\pi|y_{<t}^\pi, x) \right] \cdot R(y_{<T+1}^\pi) \\ &= -\mathbb{E}_{y_1^\pi \sim P_\theta(y_1^\pi|x)} \mathbb{E}_{y_2^\pi \sim P_\theta(y_2^\pi|y_1^\pi, x)} \dots \\ &\quad \mathbb{E}_{y_T^\pi \sim P_\theta(y_T^\pi|y_{<T}^\pi, x)} \left[ \sum_{t=1}^T \nabla_\theta \log P_\theta(y_t^\pi|y_{<t}^\pi, x) \right] \cdot R(y_{<T+1}^\pi) \\ &= -\mathbb{E}_{y_{<T+1}^\pi \sim P_\theta(y_{<T+1}^\pi|x)} \left[ \sum_{t=1}^T \nabla_\theta \log P_\theta(y_t^\pi|y_{<t}^\pi, x) \right] \cdot R(y_{<T+1}^\pi) \end{aligned} \quad (67)$$

The reward  $R(y_{<T+1}^\pi)$  will be back-propagated to every node of the computational graph via the above gradient estimator. With the gradient, the model parameters are updated by

$$\theta \leftarrow \theta + \alpha \nabla_\theta L_\theta \quad (68)$$

where  $\alpha$  is the learning rate.

As it can be seen from Eq. (67), computing gradient requires us to sample all sequences, which is not practical due to the presence of  $|\mathcal{V}|^T$  possible number of sequences. Instead, the REINFORCE algorithm approximates the expectation with a single sample, thus, the gradient is expressed as follows:

$$\nabla_\theta L_\theta \approx - \sum_{t=1}^T \nabla_\theta \log P_\theta(y_t^\pi|y_{<t}^\pi, x) \cdot R(y_{<T+1}^\pi) \quad (69)$$

One of the problems associated with this method is high variance of gradient estimator, because it makes use of only one sample to train the model. A practical solution to alleviate this problem is introducing a baseline reward [125, 104, 53] denoted by  $b$  to the gradient, i.e.,

$$\begin{aligned} \nabla_\theta L_\theta &= - \mathbb{E}_{y_{<T+1}^\pi \sim P_\theta(y_{<T+1}^\pi|x)} \left[ \sum_{t=1}^T \nabla_\theta \log P_\theta(y_t^\pi|y_{<t}^\pi, x) \right] (R(y_{<T+1}^\pi) - b) \end{aligned} \quad (70)$$

The baseline  $b$  is arbitrary function but should not depend on  $y_{<T+1}^\pi$  [53, 126]. In this way, it will not change the expectation of the gradient since

$$\mathbb{E}_{y_{<T+1}^\pi \sim P_\theta(y_{<T+1}^\pi|x)} \left[ \sum_{t=1}^T \nabla_\theta \log P_\theta(y_t^\pi|y_{<t}^\pi, x) \right] b = 0 \quad (71)$$

The complete derivations of the above equation can be found in [53]. In practice, the gradient with the baseline is approximated with

$$\nabla_\theta L_\theta \approx - \sum_{t=1}^T \nabla_\theta \log P_\theta(y_t^\pi|y_{<t}^\pi, x) (R(y_{<T+1}^\pi) - b) \quad (72)$$

Better ways of sampling a sequence and different approaches to calculate the baseline can be found in [47, 104, 127, 46].

2) *MIXER algorithm* [46]: Training seq2seq models using REINFORCE algorithm may suffer from slow convergence and can also fail due to the large action space and poor initialization (which refers to randomly initialize parameters and start with random policy). To alleviate this problem, Ranzato *et al.* [46] modified REINFORCE algorithm by incorporating the idea of curriculum learning strategy and proposed a MIXER algorithm. In this algorithm, they first trained a seq2seq model for  $N$ -epochs to ensure RL starts with a better policy. Afterwards, in each batch and for each sequence, they used the cross entropy loss for the first  $T - \Delta$  steps and REINFORCE for the remaining  $\Delta$  steps, where  $\Delta$  is an integer number. Training was continued for another  $N$ -epochs, where  $N$  is also an integer number. Then, they increased REINFORCE steps to  $2\Delta$  and continued training for another  $N$ -epochs. This process will repeat until the whole sequence is trained by REINFORCE. This algorithm has shown a better performance for greedy generation compared to XENT, DAD and E2E in the task of abstractive text summarization.

3) *Self-Critic Sequence Training (SCST)* [17, 18, 57, 53]: The main idea of SCST is to use testing time inference algorithm as the baseline function in the REINFORCE algorithm. Suppose the greedy search (see Fig. 8) is used to sample actions during testing. Then, at each training iteration, the

model generates two action sequences, in which the first one  $y_{<T+1}^{\pi, \text{greedy}}$  is from greedy search while the second one  $y_{<T+1}^{\pi}$  is sampled from a distribution  $P_{\theta}(y_{<T+1}^{\pi}|x)$ . According to SCST, baseline  $b$  is defined as reward  $R(y_{<T+1}^{\pi, \text{greedy}})$  to the first sequence. Therefore, the gradient of the loss function in SCST is expressed as

$$\nabla_{\theta} L_{\theta} \approx - \sum_{t=1}^T \nabla_{\theta} \log P_{\theta}(y_t^{\pi}|y_{<t}^{\pi}, x) (R(y_{<T+1}^{\pi}) - R(y_{<T+1}^{\pi, \text{greedy}})) \quad (73)$$

according to Eq. (72). The SCST has shown low variance and can be effectively optimized with mini-batch SGD compared to REINFORCE [53]. It has also been demonstrated to be effective in improving the performance of seq2seq models for the task of abstractive text summarization [17]. In this work, the authors used the following RL loss to train their model.

$$L_{\text{RL}} \approx - \sum_{t=1}^T \log P_{\theta}(y_t^{\pi}|y_{<t}^{\pi}, x) (R(y_{<T+1}^{\pi}) - R(y_{<T+1}^{\pi, \text{greedy}})) \quad (74)$$

Although the model performs better than those trained with XENT in terms of ROUGE scores, human-readability of generated summaries is low. To alleviate this problem, the authors also defined a mixed loss function of RL and XENT, i.e.,  $L_{\text{MIXED}} = \gamma L_{\text{RL}} + (1 - \gamma) L_{\text{XENT}}$ , where  $\gamma \in (0, 1)$  is a hyper-parameter. The model trained with the mixed loss can achieve better human-readability and ROUGE scores are still better than those obtained with XENT. They also used scheduled sampling to reducing exposure bias, in which the scheduling function is a constant ( $p_{\text{dad}} = 0.75$ ).

We have reviewed different RNN encoder-decoder architectures and training strategies in the last two sections. Now, we are at the position to generate summaries for given source articles with trained models.

#### IV. SUMMARY GENERATION

Generally speaking, the goal of summary generation is to find an optimal sequence  $y_{<T+1}^*$  such that

$$\begin{aligned} y_{<T+1}^* &= \arg \max_{y_{<T+1} \in \mathcal{Y}} \log P_{\theta}(y_{<T+1}|x) \\ &= \arg \max_{y_{<T+1} \in \mathcal{Y}} \sum_{t=1}^T \log P_{\theta}(y_t|y_{<t}, x) \end{aligned} \quad (75)$$

where  $\mathcal{Y}$  represents a set that contains all possible sequences (summaries). However, since it has  $|\mathcal{Y}|^T$  elements, the exact inference is intractable in practice [15]. Here,  $\mathcal{Y}$  represents the output vocabulary. In this section, we review the beam search algorithm and its extensions for approximating the exact inference.

##### A. Greedy and Beam Search

As shown in Fig. 8, we can generate a sub-optimal sequence with greedy search, i.e.,

$$y_t^* = \arg \max_{y_t \in \mathcal{V}} \log P_{\theta}(y_t|y_{<t}, x) \quad (76)$$

at each decoding step  $t$ . Although greedy search is computationally efficient, human-readability of generated summaries is low.

Beam search algorithm is a compromise between greedy search and exact inference and has been commonly employed in different language generation tasks [15, 21, 12]. Beam search is a graph-search algorithm that generates sequences from left to right by retaining only  $B$  top scoring (top- $B$ ) sequence-fragments at each decoding step. More formally, we denote decoded top- $B$  sequence fragments, also known as hypotheses [15], at time-step  $t - 1$  as  $y_{<t,1}, y_{<t,2}, \dots, y_{<t,B}$  and their scores as  $S_{<t,1}^{\text{bm}}, S_{<t,2}^{\text{bm}}, \dots, S_{<t,B}^{\text{bm}}$ . For each fragment  $y_{<t,b}$ , we first calculate  $P_{\theta}(y_{t,b}^{\text{cand}}|y_{<t,b}, x)$ , which determines  $B$  most probable words  $y_{t,b,1}^{\text{cand}}, y_{t,b,2}^{\text{cand}}, \dots, y_{t,b,B}^{\text{cand}}$  to expand it. The score for each expanded fragment, i.e., new hypotheses,  $y_{<t+1,b,b'}^{\text{cand}}$  can then be updated with either

$$S_{t,b,b'}^{\text{cand}} = S_{<t,b}^{\text{bm}} \times P_{\theta}(y_{t,b,b'}^{\text{cand}}|y_{<t,b}, x) \quad (77)$$

where  $S_{<t,b}^{\text{bm}}$  is initialized with 1, or

$$S_{t,b,b'}^{\text{cand}} = S_{<t,b}^{\text{bm}} + \log P_{\theta}(y_{t,b,b'}^{\text{cand}}|y_{<t,b}, x) \quad (78)$$

where  $S_{<t,b}^{\text{bm}}$  is initialized with 0. Here,  $b$  and  $b'$  are labels of a current hypothesis and a word candidate, respectively. This yields  $B \times B$  expanded fragments, i.e., new hypotheses, in which only the top- $B$  of them along with their scores are retained for the next decoding step. This procedure will be repeated until ‘EOS’ token is generated. In Algorithm 1, we show pseudocodes of a beam search algorithm for generating summaries with attention-based seq2seq models given the beam size of  $B$  and batch size of 1.

##### B. Diverse Beam Decoding

Despite widespread applications, beam search algorithm suffered from lacking of diversity within a beam [89, 93, 94]. In other words, the top- $B$  hypotheses may differ by just a couple tokens at the end of sequences, which not only affects the quality of generated sequences but also wastes computational resources [90, 94]. In this section, we briefly introduce some studies that aim to improve decoding by increasing the diversity of the beam search algorithm.

###### 1) Maximum Mutual Information (MMI) [90, 92, 93]:

The MMI based methods were originally proposed for neural conversation models and then applied to other tasks, such as machine translation and summarization [92, 93]. The basic intuition here is that a desired model should not only take into account the dependency of a target on a source, but also should consider the likelihood of the source for a given target, which is achieved by replacing the log-likelihood of the target, i.e.,  $\log P_{\theta}(y|x)$  in Eq. (75), with pairwise mutual information of the source and target, defined by  $\log \frac{P_{\theta}(y,x)}{P_{\theta}(x)P_{\theta}(y)}$ . During the training, model parameters are learned by maximizing mutual information. When generating sequences, the objective is expressed as follows:

$$\begin{aligned} y^* &= \arg \max_{y \in \mathcal{Y}} \log \frac{P_{\theta}(y, x)}{P_{\theta}(x)P_{\theta}(y)} \\ &= \arg \max_{y \in \mathcal{Y}} (\log P_{\theta}(y|x) - \log P_{\theta}(y)) \end{aligned} \quad (79)$$

---

**Algorithm 1:** Beam search algorithm for decoding the basic attention-based seq2seq models.

---

**Input:** Source article  $x$ , beam size  $B$ , summary length  $T$ , model parameters  $\theta$ ;

**Output:**  $B$ -best summaries;

```

1 Initialize:
2 Output sequences  $Q^{\text{seq}} = [\text{SOS}]_{B \times T}$ ;
3 Accumulated probabilities  $Q^{\text{prob}} = [1.0]_{B \times 1}$ ;
4 The last decoded tokens  $Q^{\text{word}} = [\text{SOS}]_{B \times 1}$ ;
5 States (hidden and cell states for LSTM)  $Q^{\text{states}} = [0.0]_{B \times |h_t^d|}$ ;
6 Context vectors  $Q^{\text{ctx}} = [0.0]_{B \times |z_t^e|}$ ;

7 Compute  $(h_1^e, h_2^e, \dots, h_J^e)$  with encoder;
8 Update  $Q^{\text{states}}$  with encoder states;
9 for  $t=1, T$  do
10   Initialize candidates  $Q^{\text{cand,seq}}, Q^{\text{cand,prob}}, Q^{\text{cand,word}}, Q^{\text{cand,states}}, Q^{\text{cand,ctx}}$  by repeating  $Q^{\text{seq}}, Q^{\text{prob}}, Q^{\text{word}}, Q^{\text{states}}$  and
       $Q^{\text{ctx}}$   $B$  times, respectively;
11   for  $b=1, B$  do
12     Compute  $P_\theta(y_{t,b}^{\text{cand}} | y_{<t,b}, x)$  using decoder LSTM cell with input  $(h_1^e, h_2^e, \dots, h_J^e)$ ,  $Q_b^{\text{word}}, Q_b^{\text{states}}$  and  $Q_b^{\text{ctx}}$ ;
13     Select the top- $B$  candidate words  $y_{t,b,b'}^{\text{cand}}$ , where  $b' = 1, 2, \dots, B$ ;
14     Select corresponding probability  $P_\theta(y_{t,b,b'}^{\text{cand}} | y_{<t,b}, x)$ , hidden states  $h_{t,b,b'}^d$ , cell states  $c_{t,b,b'}^d$  and context vector
         $z_{t,b,b'}^e$ ;
15     Update elements of  $Q_{b',b}^{\text{cand,seq}}, Q_{b',b}^{\text{cand,word}}$  with  $y_{t,b,b'}^{\text{cand}}$ ;
16     Update elements of  $Q_{b',b}^{\text{cand,states}}$  with  $h_{t,b,b'}^d$  and  $c_{t,b,b'}^d$ ;
17     Update elements of  $Q_{b',b}^{\text{cand,ctx}}$  with  $z_{t,b,b'}^e$ ;
18     Update  $Q_{b',b}^{\text{cand,prob}}$  with Eq.(77);
19   end
20   Flatten  $Q^{\text{cand,prob}}$  and choose  $B$  best hypotheses;
21   Update  $Q_t^{\text{seq}}, Q^{\text{prob}}, Q^{\text{word}}, Q^{\text{states}}, Q^{\text{ctx}}$  with corresponding candidates.
22 end

```

---

However, it is obvious that calculating  $P_\theta(y)$  is intractable. Thus, several approximation methods have been proposed in the literature to alleviate this problem.

**MMI-antiLM [90]:** Eq. (79) can be further generalized as

$$y^* = \arg \max_{y \in \mathcal{Y}} (\log P_\theta(y|x) - \lambda \log P_\theta(y)) \quad (80)$$

where parameter  $\lambda$  controls the influence of the penalty term.  $P_\theta(y)$  is expanded as follows:

$$P_\theta(y) = \prod_{t=1}^T P_\theta(y_t | y_{<t}) g(t) \quad (81)$$

where  $g(t)$  is a monotonically non-increasing weight, which encourages diversity of first few words of sequences. Intuitively, these words play a significant role in determining the remainder of a sequence and penalizing them leads to more diverse results. As the influence of the generated words to the diversity declines, the model relies more on the language model due to syntactic and grammatical consideration. In this approach,  $g(t)$  is defined as  $g(t) = \begin{cases} 1 & t \leq \kappa \\ 0 & t > \kappa \end{cases}$  and

$0 < \kappa < T$  is an integer threshold. This is the first method to approximate this objective.

**MMI-bidi [90]:** According to the Baye's theorem, we have

$$\log P_\theta(y) = \log P_\theta(y|x) + \log P_\theta(x) - \log P_\theta(x|y) \quad (82)$$

Substituting Eq. (82) into Eq. (80) yields

$$y^* = \arg \max_{y \in \mathcal{Y}} ((1 - \lambda) \log P_\theta(y|x) + \lambda \log P_\theta(x|y)) \quad (83)$$

which can be viewed as a trade-off between the dependency of source-on-target and target-on-source. However, it is intractable to generate sequences directly using Eq. (83), because calculating  $P_\theta(x|y)$  requires that the generation is completed. An approximation method has been proposed to remedy this problem by first constructing the so-called  $N$ -best lists with the first term of Eq. (83), and then re-ranking the list by jointly taking into account both  $\log P_\theta(y|x)$  and  $\log P_\theta(x|y)$ .

**Mutual Information via Re-ranking (MIR) [90, 92]:** The above approximation method can be summarized by the following three steps:

- Train two seq2seq models, one for  $P_{\theta_1}(y|x)$  and the other for  $P_{\theta_2}(x|y)$ , where  $\theta_1$  and  $\theta_2$  are the model parameters.
- Generate a diverse  $N$ -best list of sequences based on  $P_{\theta_1}(y|x)$ . To achieve this goal, the method for calculating the scores for beam search algorithm has been modified as

$$S_{k,k'}^{\text{beam}} = b_{<t,k} + \log P_\theta(y_{t,k'} | y_{<t,k}, x) - \gamma k' \quad (84)$$

By adding the last term  $\gamma k'$ , the model explicitly encourages hypotheses from different parents, i.e., different  $k$ , which results in more diverse results. Therefore, parameter  $\gamma$  is

also known as the diversity rate which indicates the degree of diversity integrated into beam search algorithm [93].

- Re-rank  $N$ -best list by linearly combining  $P_{\theta_1}(y|x)$  and  $P_{\theta_2}(x|y)$ . The ranking score for each candidate in the  $N$ -best list is defined as follows:

$$S_{\text{rank}}(y) = \log P_{\theta_1}(y|x) + \lambda \log P_{\theta_2}(x|y) + \beta \Omega(y) \quad (85)$$

where  $\Omega(y)$  is a task-specific auxiliary term.  $\lambda$  and  $\beta$  are parameters that can be learned using minimum error rate training [128] on the development dataset.

**Automatically Learning Diversity Rate [93]:** The modified score function (Eq. (84)) for beam search algorithm enables us to generate diverse  $N$ -best lists. However, the diversity rate  $\gamma$  is hard-coded into the algorithm, which yields sub-optimal results. To alleviate this problem, a RL-based method has been proposed to find the best diversity rate  $\gamma$  with respect to the given evaluation metrics. In this model, the diversity rate is defined as a function of input  $x$ , i.e.,  $\gamma = \gamma(x)$ . It can be only chosen from a finite list of real numbers, denoted by  $\Gamma^{14}$ , by a policy network  $\pi(\gamma(x) = \gamma'|x)$ , which is the distribution over  $|\Gamma|$  classes and defined as follows:

$$\pi(\gamma(x) = \gamma'|x) = \frac{\exp(h_x^\top \cdot h_{\gamma'})}{\sum_{j=1}^{|\Gamma|} \exp(h_x^\top \cdot h_{\Gamma_j})} \quad (86)$$

where  $h_x$  is a vector representation of given input  $x$  that can be calculated by Eq. 1. After an action is taken, we can follow step (2) and (3) of MIR to select the best output, which is then used to calculate the reward  $R(\gamma(x) = \gamma')$  (e.g., BLEU or ROUGE scores). Finally, the optimal diversity rate policy is obtained by maximizing the expectation of final reward using the REINFORCE algorithm [51]. This simple method enables the model to adjust its diversity rate in response to input and generate better sequences than the basic diverse decoding approach.

2) *Diverse Beam Search (DBS) [94, 95]:* DBS is another approach that aims to increase the diversity of standard beam search algorithm. It first partitions the hypotheses into  $G$  groups. Then, at each decoding step, it sequentially performs a beam search on each group based on a dissimilarity augmented scoring function

$$S_{t,b,b'}^{\text{cand}} = S_{<t,b}^{\text{cand}} + \log P_{\theta}(y_{t,b,b'}^{\text{cand},g} | y_{<t,b}^g, x) + \lambda_g \Delta(y_{<t+1,b,b'}^{\text{cand},g}; y_{<t+1}^1, \dots, y_{<t+1}^{g-1}) \quad (87)$$

where  $\lambda_g \geq 0$  is a parameter.  $\Delta(y_{<t+1,b,b'}^{\text{cand},g}; y_{<t+1}^1, \dots, y_{<t+1}^{g-1})$  represents a diversity function, which measures the dissimilarity or distance between candidate  $y_{<t+1,b,b'}^{\text{cand},g}$  in group  $g$  and sequences in groups from 1 to  $g-1$ . Standard beam search is applied to group 1. Intuitively, the generated sequences in different groups are very different from each other due to the penalty of diversity functions. In [95], DBS has been combined with pointer-generator network [12] to improve the diversity of the model produced summaries.

## V. CONVOLUTIONAL SEQ2SEQ MODELS

RNN encoder-decoder models have been commonly used in sequence modeling and language generation tasks. However, they are limited to model relatively short sequences due to their intrinsic sequential dependency in the time-step dimension (see Fig. 2). In other words, for longer sequences, they require more GPU memory to store computational graphs and more training time. Other network architectures, such as feed-forward network and CNN, also have similar problems, but these problems can be alleviated by parallelizing these algorithms along the time-step dimension [67, 63]. In this section, we briefly review the recently proposed convolutional seq2seq model (ConvS2S) [67], which have also shown good performance on the task of abstractive text summarization [67, 72, 119]. It should be noted that most notations in this section are self-consistent and independent from other sections.

### A. Position Embedding

We still adopt the same notations as describe in Section II-A for the source article  $x = (x_1, x_2, \dots, x_{|x|})$  and target summary  $y = (y_1, y_2, \dots, y_{|y|})$ , where  $x_j$  and  $y_t$  are one-hot representations of tokens, and  $|x|$  and  $|y|$  represent lengths of sequences. The word embedding of each token  $x_j$  is denoted by  $E_{x_j}^{\text{wd}}$ . In ConvS2S, the embedding layer also incorporates embeddings of absolute position of tokens in sequences, i.e.,  $E_{x_j}^{\text{pos}}$ , to keep track of order information. Therefore, the final embedding is a combination of word and position embeddings, i.e.,  $E_{x_j} = E_{x_j}^{\text{wd}} + E_{x_j}^{\text{pos}}$ . Embeddings for tokens  $E_{y_t}$  in summaries can be obtained in the same manner.

### B. Convolutional Encoder and Decoder

In ConvS2S architecture, both encoder and decoder are constructed with a multi-layer CNN, where each layer has a one-dimensional convolution followed by some non-linearity. More specifically, for layer  $l$ , a convolution is a linear operation that takes a vector  $v^l \in \mathbb{R}^{kd}$  as input and outputs another vector  $g^l \in \mathbb{R}^{2d}$ , where  $d$  is the dimension of embeddings and  $k$  is the size of the convolutional kernel. Formally, it can be expressed as follows:

$$g^l = W_{\text{conv}}^l v^l + b_{\text{conv}}^l \quad (88)$$

where  $W_{\text{conv}}^l$  and  $b_{\text{conv}}^l$  are learnable parameters.  $v^l$  is the concatenation of  $k$  rows of input matrix  $H^l \in \mathbb{R}^{|x| \times B}$ <sup>15</sup>, i.e.,  $v^l = [h_{j-k/2}^l, \dots, h_{j+k/2}^l]$ , and  $h_j^l$  is the  $j$ -th row of  $H^l$ . We first transform output vector  $g^l \in \mathbb{R}^{2d}$  to a matrix  $G^l \in \mathbb{R}^{2 \times d}$ . Then, a gated linear unit (GLU) [73] is applied on top of  $G^l$  as non-linearity, in which the gate mechanism is defined as follows:

$$h_j^{l+1} = G_1^l \otimes \sigma(G_2^l) \quad (89)$$

where  $\otimes$  is an element-wise multiplication operator, and  $\sigma(G_2^l)$  is a sigmoid function which controls the information flow of the current context. To achieve better performance, residual connections [129] have also been introduced in each layer from input to output, i.e.,

$$h_j^{l+1} = G_1^l \otimes \sigma(G_2^l) + h_j^l \quad (90)$$

<sup>15</sup>When  $l = 1$ ,  $H^l$  is the embedding matrix. When  $l > 1$ ,  $H^l$  is output of layer  $l-1$ . Zero padding has to be applied to  $H^l$ .

<sup>14</sup>‘chosen’ and  $\Gamma$  can be viewed as an action and an action space in RL.

### C. Multi-Step Attention

Another novel element in ConvS2S is multi-step attention mechanism, which has been applied to each layer of the decoder. Suppose we denote the output of the last layer  $u$  of encoder network by  $h^u = (h_1^u, \dots, h_{|x|}^u)$  (encoder hidden states) and output of layer  $l$  of decoder network by  $d^l = (d_1^l, \dots, d_{|y|}^l)$  (decoder hidden states). To compute the attention weights, we first combine current decoder state  $d_t^l$  with embedding of the previous token  $E_{y_{t-1}}$  as

$$z_t^l = W_{\text{dec}}^l d_t^l + b_{\text{dec}}^l + E_{y_{t-1}} \quad (91)$$

where weight matrix  $W_{\text{dec}}^l$  and bias  $b_{\text{dec}}^l$  are parameters. Then, the attention weights can be obtained by

$$\alpha_{tj}^l = \frac{\exp(z_t^l \cdot h_j^u)}{\sum_{k=1}^{|x|} \exp(z_t^l \cdot h_k^u)} \quad (92)$$

Finally, we can take the weighted sum of encoder hidden states and embeddings to calculate conditional input  $c_t^l$ <sup>16</sup>

$$c_t^l = \sum_{j=1}^{|x|} \alpha_{tj}^l (h_j^u + E_{x_j}). \quad (93)$$

Conditional input vector  $c^l = (c_1^l, \dots, c_{|y|}^l)$  will be further added to the corresponding output  $d^l = (d_1^l, \dots, d_{|y|}^l)$  to compute  $d^{l+1} = (d_1^{l+1}, \dots, d_{|y|}^{l+1})$  with Eqs. (88) and (90). At each time step, the model performs multiple attention ‘hops’ in contrast with single step attention [23, 21], and by applying convolution (i.e., Eq. (88)), the decoder can also take into account the previously attended inputs.

## VI. IMPLEMENTATIONS AND EXPERIMENTS

Apart from a comprehensive literature survey and a detailed review of different techniques for network structures, training strategies and summary generations, we have also developed an open-source library, namely, NATS<sup>17</sup>, based on RNN seq2seq framework for abstractive text summarization. In this section, we first introduce the details of our implementations and then systematically experiment with different network elements and hyper-parameters on three public available datasets, i.e., CNN/Daily Mail, Newsroom, and Bytecup.

### A. Implementations

The NATS is equipped with following important features<sup>18</sup>:

- **Attention based seq2seq framework**. We implemented the attention based seq2seq model shown in Fig. 3. Encoder and decoder can be chosen to be either LSTM or GRU. The attention scores can be calculated with one of three alignment methods given in Eq. (6).
- **Pointer-generator network\***. Based on the attention based seq2seq framework, we implemented pointer-generator network discussed in Section II-C4.
- **Intra-temporal attention mechanism\***. The temporal attention can work with all three alignment methods.

<sup>16</sup>It is similar to the context vector in RNN seq2seq models.

<sup>17</sup><https://github.com/tshi04/NATS>

<sup>18</sup>\* indicates users can switch on/off a feature or have multiple options.

- **Intra-decoder attention mechanism\***. The alignment method for intra-decoder attention is the same as that for the attention mechanism.
- **Coverage mechanism\***. To handle the repetition problem, we also implemented the coverage mechanism discussed in Section II-D3. If coverage is switched off, the coverage loss will be set to 0.
- **Weight sharing mechanism\***. As discussed in Section II-F1, weight sharing mechanism can boost the performance using significantly fewer parameters.
- **Beam search algorithm**. We implemented an efficient beam search algorithm that can also handle the case when the batch size  $> 1$ .
- **Unknown words replacement\***. Similar to [107], we implemented a heuristic unknown words replacement technique to boost the performance. Theoretically, a pointer-generator network may generate OOV words even with the copying mechanism, because  $\langle \text{unk} \rangle$  is still in the extended vocabulary. Thus, after the decoding is completed, we manually check  $\langle \text{unk} \rangle$  in summaries and replace them with words in source articles using attention weights. This meta-algorithm can be used for any attention-based seq2seq model.

### B. Datasets

1) *CNN/Daily Mail Dataset*: CNN/Daily Mail dataset<sup>19</sup> consists of more than 300K news articles and each of them is paired with several highlights, known as multi-sentence summaries [14, 12]. We have summarized the basic statistics of the dataset in Table III. There are primarily two versions of this dataset. The first version anonymizes name entities [14], while the second one keeps the original texts [12]. In this paper, we used the second version and obtained processed data from See *et al.* [12]<sup>20</sup>.

TABLE III: Basic statistics of the CNN/Daily Mail dataset.

	Train	Validation	Test
# pairs	287,227	13,368	11,490
Article Length	751	769	778
Summary Length	55	61	58

2) *Newsroom Dataset*: The Cornell Newsroom dataset<sup>21</sup> [100] was recently released and consists of 1.3 million article-summary pairs, out of which 1.2 million of them are publicly available for training and evaluating summarization systems. We first used newsroom library<sup>22</sup> to scrape and extract the raw data. Then, texts were tokenized with NLTK<sup>23</sup>, SpaCy<sup>24</sup> and Stanford CoreNLP<sup>25</sup> packages, respectively. We developed a data processing tool that is make available at <sup>26</sup> to tokenize texts and prepare input for

<sup>19</sup><https://github.com/abisee/cnn-dailymail>

<sup>20</sup><https://github.com/JafferWilson/Process-Data-of-CNN-DailyMail>

<sup>21</sup><https://summari.es/>

<sup>22</sup><https://github.com/clc-lab/newsroom>

<sup>23</sup><https://www.nltk.org/>

<sup>24</sup><https://spacy.io/>

<sup>25</sup><https://stanfordnlp.github.io/CoreNLP/>

<sup>26</sup>[https://github.com/tshi04/NATS/tree/master/tools/newsroom\\_process](https://github.com/tshi04/NATS/tree/master/tools/newsroom_process)



NATS. In this survey, we created six datasets with different tokenizers, three of which are for text summarization and the rest of them are used for headline generation. The basic statistics of these datasets are shown in Table IV.

TABLE IV: Basic statistics of the Newsroom dataset.

Tokenizer	Characteristics	Train	Validation	Test
CoreNLP	# pairs	992,985	108,612	108,655
	Article Length	756.25	750.40	749.20
	Summary Length	29.99	30.35	30.25
	Headline Length	9.70	9.69	9.69
NLTK	Article Length	756.59	749.78	748.89
	Summary Length	30.01	30.36	30.27
	Headline Length	9.63	9.62	9.63
SpaCy	Article Length	773	766	765
	Summary Length	30.65	31.02	30.92
	Headline Length	9.92	9.91	9.91

3) *Bytecup Dataset*: Byte Cup 2018 International Machine Learning Contest<sup>27</sup> released a new dataset (which will be referred to as Bytecup dataset in this survey) for the headline/title generation task. It consists of 1.3 million pieces of articles, out of which 1.1 million are released for training. In our experiments, we create training, development and testing sets (0.8/0.1/0.1) based on this training dataset. Texts are tokenized using Stanford CoreNLP package and prepared with our data processing tool<sup>28</sup>. The basic statistics of the dataset are shown in Table V.

TABLE V: Basic statistics of the Bytecup dataset.

Characteristics	Train	Validation	Test
# pairs	892,734	111,592	111,592
Article Length	639.91	640.12	638.82
Summary Length	12.02	12.03	12.04

### C. Parameter Settings

In all our experiments, we set the dimension of word embeddings and hidden states (for both encoder and decoder) as 128 and 256, respectively. During training, the embeddings are learned from scratch. Adam [130] with hyper-parameter  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and  $\epsilon = 10^{-8}$  is used for stochastic optimization. Learning rate is fixed to 0.0001 and mini-batches of size 16 are used. Gradient clipping is also used with a maximum gradient norm of 2.0. For all datasets, the vocabulary consists of 50K words and is shared between source and target. For the CNN/Daily Mail dataset, we truncate source articles to 400 tokens and limit the length of summaries to 100 tokens. For the Newsroom dataset, source articles, summaries and headlines are truncated to 400, 50 and 20, respectively. For the Bytecup dataset, lengths of source articles and headlines are also limited to 400 and 20 tokens, respectively. During training, we run 35 epochs for the CNN/Daily Mail dataset and 20 epochs for the Newsroom and Bytecup dataset. During testing, we set the size of a beam to 5.

### D. ROUGE Evaluations

Recall-Oriented Understudy for Gisting Evaluation (ROUGE) scores were first introduced in [16] and have become standard metrics for evaluating abstractive text summarization models. They determine the quality of summarization by counting the number of overlapping units (i.e.,  $n$ -grams, word sequences and word pairs) between machine generated and golden-standard (human-written) summaries [16]. Within all different ROUGE measures, ROUGE-1 (unigram), ROUGE-2 (bigram) and ROUGE-L (longest common subsequence) have been most widely used for single-document abstractive summarization [12]. In this paper, different models are evaluated using pyrouge<sup>29</sup> package, which provides precision, recall and F-score for these measures.

### E. Experiments on CNN/Daily Mail dataset

In the past few years, the CNN/Daily Mail dataset has become a standard benchmark dataset used for evaluating the performance of different summarization models that can generate multi-sentence summaries for relatively longer documents [12, 81, 131, 14, 17]. Hence, before discussing our experimental results, we first review ROUGE performance of different abstractive models on CNN/Daily Mail dataset in Table VII. In our experiments, we systematically investigated the effects of six network components in the seq2seq framework on summarization performance, including (i) alignment methods in the attention mechanism, (ii) pointing mechanism, (iii) intra-temporal attention, (iv) intra-decoder attention, (v) weight sharing and (vi) coverage mechanism.

Our experimental results are shown in Table VIII. To effectively represent different models, ID of each model consists of a letter followed by five binary-indicators, corresponding to the six important components. The letters ‘G’, ‘D’ and ‘C’ denote alignment methods ‘general’, ‘dot’ and ‘concat’, respectively. 1 and 0 indicates if a component is switched on or off, respectively. At first, it can be clearly seen that the performance of three basic attention based models (i.e., G00000, D00000, C00000) are close to each other. In these tests, we still keep the OOV tokens in generated summaries when performing the ROUGE evaluations, which results in relatively lower ROUGE precision scores. Therefore, ROUGE F-scores may be lower than those reported in the literature [14, 12]. Comparing G10000 with G00000, and C10000 with C00000, we find that pointing mechanism significantly improves the performance of attention based seq2seq models. By analyzing summaries, we observed that most of the tokens are copied from source articles, which results in summaries that are similar to the ones generated by the extractive models<sup>30</sup>. As discussed in Section II-C4, another advantage of pointing mechanism is that it can effectively handle OOV tokens. From the results of G10000 and C10000, we see that the pointing mechanism works better when used in the ‘concat’ alignment method.

<sup>29</sup><https://pypi.python.org/pypi/pyrouge/0.1.0>

<sup>30</sup>The extractive models attempt to extract sentences from the source articles.

<sup>27</sup><https://www.biendata.com/competition/bytecup2018/>

<sup>28</sup>[https://github.com/tshi04/NATS/tree/master/tools/bytecup\\_process](https://github.com/tshi04/NATS/tree/master/tools/bytecup_process)

TABLE VI: ROUGE scores of different models on the entity-anonymized CNN/Daily Mail dataset.

Year	Reference	Model	ROUGE-1	ROUGE-2	ROUGE-L
2016	Nallapati <i>et al.</i> [14]	words-lvt2k-temp-att	35.46	13.30	32.65
2017	Paulus <i>et al.</i> [17]	Intra-attention + RL	41.16	15.75	39.08
		Intra-attention + RL + ML	39.87	15.82	36.90
	Tan <i>et al.</i> [75]	Graph-based attention	38.01	13.90	34.00
	Liu <i>et al.</i> [117]	Adversarial Training	39.92	17.65	36.71
2018	Li <i>et al.</i> [82]	KIGN + Prediction-guide	38.95	17.12	35.68
	Chen <i>et al.</i> [59]	RNN-ext + abs + RL + rerank	39.66	15.85	37.34
	Kryscinski <i>et al.</i> [120]	ML + RL ROUGE + Novel, with LM	40.02	15.53	37.44

TABLE VII: ROUGE scores of different models on the non-anonymized CNN/Daily Mail dataset.

Year	Reference	Model	ROUGE-1	ROUGE-2	ROUGE-L
2017	See <i>et al.</i> [12]	Pointer-generator	36.44	15.66	33.42
		Pointer-generator + coverage	39.53	17.28	36.38
	Fan <i>et al.</i> [72]	ConvS2S	39.75	17.29	36.54
2018	Celikyilmaz <i>et al.</i> [18]	Deep Communicating Agents (m7)	41.69	19.47	37.92
	Chen <i>et al.</i> [59]	RNN-ext + abs + RL + rerank	40.88	17.80	38.54
	Hsu <i>et al.</i> [80]	End2end w/ inconsistency loss	40.68	17.97	37.13
	Pasunuru <i>et al.</i> [57]	ROUGESal + Ent (RL)	40.43	18.00	37.10
	Guo <i>et al.</i> [87]	Two-Layer Baseline (Pointer+Coverage) + Entailment Gen. + Question Gen.	39.81	17.64	36.54
	Cibils <i>et al.</i> [95]	Pointer-Generator Network with DBS + EmbedRank decoder	40.19	17.09	36.63
	Gehrmann <i>et al.</i> [43]	Bottom-Up Summarization	41.22	18.68	38.34
	Jiang <i>et al.</i> [42]	RL + pg + cbdec	40.66	17.87	37.06
	Chung <i>et al.</i> [122]	Main Pointer Generator + Coverage	40.12	17.74	36.82
	Kryscinski <i>et al.</i> [120]	ML + RL ROUGE + Novel, with LM	40.19	17.38	37.52

The remaining four components are tested upon the pointer-generator network. By comparing G11000 and G10000, we see that the intra-temporal attention increases almost 3 ROUGE points. This might be because of its capability of reducing repetitions. However, most of our models that combine intra-temporal attention with ‘concat’ failed during training after a few epochs. Thus, we did not report these results. As to intra-decoder attention, we observe from G11000, G11010, G11100 and G11110 that it does not boost the performance of the model before adding weight sharing mechanism. However, in the case of ‘concat’, the models with intra-decoder attention have a better performance. Weight sharing mechanism does not always boost the performance of the models (according to the comparison of C10100 and C10110). However, as aforementioned, models that adopt weight sharing mechanism have much fewer parameters. Finally, we find the coverage mechanism can significantly boost performance by at least 2 ROUGE points, which is consistent with the results presented in [12].

#### F. Experiments on Newsroom and Bytecup Datasets

We also tested NATS toolkit on the Newsroom dataset, which was released recently. In our experiments, we tokenized the raw data with three different packages and generated three versions of the dataset for the task of text summarization and three versions for the task of headline generation. Experi-

mental results obtained with the models G11110 and C10110 on the released testing set [100] are shown in Table IX and Table X. It can be observed that G11110 performs better than C10110 on CNN/Daily Mail data from Table VIII, however, C10110 achieves better ROUGE scores in both text summarization and headline generation tasks on Newsroom dataset. Further more, we also evaluated the trained models on the CNN/Daily Mail testing set and the results are summarized in Table XI. It can be observed that C10110 performs better than G11110. Finally, we summarize our results for the Bytecup headline generation dataset in Table XII. C10110 still achieves slightly better ROUGE scores than G11110.

## VII. CONCLUSION

Being one of the most successful applications of seq2seq models, neural abstractive text summarization has become a prominent research topic that has gained a lot of attention from both industry and academia. In this paper, we provided a comprehensive survey on the recent advances of seq2seq models for the task of abstractive text summarization. This work primarily focuses on the challenges associated with neural network architectures, model parameter inference mechanisms and summary generation procedures, and the solutions of different models and algorithms. We also provided a taxonomy of these topics and an overview of different seq2seq models for the abstractive text summarization. As part of this survey,

TABLE VIII: ROUGE scores on the CNN/Daily Mail dataset in our experiments.

Model ID	Attention	Pointer-generator	Intra-temporal	Intra-decoder	Weight sharing	Coverage	ROUGE-1	ROUGE-2	ROUGE-L
G00000	general	-	-	-	-	-	27.62	10.15	25.81
G10000	general	✓	-	-	-	-	33.85	14.08	31.33
G11000	general	✓	✓	-	-	-	36.78	15.82	34.05
G11010	general	✓	✓	-	✓	-	37.07	15.93	34.22
G11100	general	✓	✓	✓	-	-	36.74	15.78	33.98
G11110	general	✓	✓	✓	✓	-	<b>37.60</b>	<b>16.27</b>	<b>34.77</b>
D00000	dot	-	-	-	-	-	27.60	10.19	25.81
D11100	dot	✓	✓	✓	-	-	37.00	15.87	34.24
D11110	dot	✓	✓	✓	✓	-	<b>37.65</b>	<b>16.25</b>	<b>34.85</b>
C00000	concat	-	-	-	-	-	27.61	10.16	25.78
C10000	concat	✓	-	-	-	-	35.55	15.18	32.84
C10001	concat	✓	-	-	-	✓	38.64	16.70	35.63
C10100	concat	✓	-	✓	-	-	36.51	15.75	33.70
C10101	concat	✓	-	✓	-	✓	<b>39.23</b>	<b>17.28</b>	36.02
C10110	concat	✓	-	✓	✓	-	36.46	15.68	33.69
C10111	concat	✓	-	✓	✓	✓	39.14	17.13	<b>36.04</b>

TABLE IX: ROUGE scores on the Newsroom-summary dataset.

Tokenizer	Model	ROUGE-1	ROUGE-2	ROUGE-L
CoreNLP	G11110	38.86	27.33	35.69
NLTK	G11110	38.87	27.36	35.83
SpaCy	G11110	39.11	27.54	35.99
	C10110	39.36	27.86	36.35

TABLE X: ROUGE scores on the Newsroom-headline dataset.

Tokenizer	Model	ROUGE-1	ROUGE-2	ROUGE-L
CoreNLP	G11110	25.82	11.71	23.94
NLTK	G11110	25.31	11.49	23.55
SpaCy	G11110	25.23	11.46	23.72
	C10110	26.56	12.31	25.04

we developed an open source toolkit, namely, NATS, which is equipped with several important features, including attention, pointing mechanism, repetition handling, and beam search. In our experiments, we first summarized the experimental results of different seq2seq models in the literature on the widely used CNN/Daily Mail dataset. We also conducted extensive experiments on this dataset using NATS to examine the effectiveness of different neural network components. Finally, we established benchmarks for two recently released datasets, i.e., Newsroom and Bytecup.

#### ACKNOWLEDGMENTS

This work was supported in part by the US National Science Foundation grants IIS-1619028, IIS-1707498 and IIS-1838730.

#### REFERENCES

[1] D. R. Radev, E. Hovy, and K. McKeown, “Introduction to the special issue on summarization,” *Computational linguistics*, vol. 28, no. 4, pp. 399–408, 2002.

TABLE XI: ROUGE scores on the CNN/Daily Mail dataset.

Tokenizer	Model	ROUGE-1	ROUGE-2	ROUGE-L
CoreNLP	G11110	32.92	12.34	28.58
NLTK	G11110	33.86	12.85	29.51
SpaCy	G11110	32.17	11.85	27.78
	C10110	34.28	13.11	29.94

TABLE XII: ROUGE scores on the Bytecup dataset.

Model	ROUGE-1	ROUGE-2	ROUGE-L
G11110	39.04	22.72	35.90
C10110	39.13	22.98	36.14

- [2] M. Allahyari, S. Pouriyeh, M. Assefi, S. Safaei, E. D. Trippe, J. B. Gutierrez, and K. Kochut, “Text summarization techniques: A brief survey,” *arXiv preprint arXiv:1707.02268*, 2017.
- [3] I. Mani and M. T. Maybury, *Advances in automatic text summarization*. MIT press, 1999.
- [4] M. Gambhir and V. Gupta, “Recent automatic text summarization techniques: a survey,” *Artificial Intelligence Review*, vol. 47, no. 1, pp. 1–66, 2017.
- [5] R. M. Verma and D. Lee, “Extractive summarization: Limits, compression, generalized model and heuristics,” *Computación y Sistemas*, vol. 21, 2017.
- [6] N. Bhatia and A. Jaiswal, “Automatic text summarization and it’s methods-a review,” in *Cloud System and Big Data Engineering (Confluence)*, 2016 6th International Conference. IEEE, 2016, pp. 65–72.
- [7] E. Lloret and M. Palomar, “Text summarisation in progress: a literature review,” *Artificial Intelligence Review*, vol. 37, no. 1, pp. 1–41, 2012.
- [8] H. Saggion and T. Poibeau, “Automatic text summarization: Past, present and future,” in *Multi-source, multilingual information extraction and summarization*. Springer, 2013, pp. 3–21.
- [9] A. Nenkova, K. McKeown *et al.*, “Automatic sum-

- marization,” *Foundations and Trends® in Information Retrieval*, vol. 5, no. 2–3, pp. 103–233, 2011.
- [10] D. Das and A. F. Martins, “A survey on automatic text summarization,” *Literature Survey for the Language and Statistics II course at CMU*, vol. 4, pp. 192–195, 2007.
- [11] Y. Wu and B. Hu, “Learning to extract coherent summary via deep reinforcement learning,” 2018.
- [12] A. See, P. J. Liu, and C. D. Manning, “Get to the point: Summarization with pointer-generator networks,” in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 2017, pp. 1073–1083.
- [13] Q. Zhou, N. Yang, F. Wei, S. Huang, M. Zhou, and T. Zhao, “Neural document summarization by jointly learning to score and select sentences,” in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, vol. 1, 2018, pp. 654–663.
- [14] R. Nallapati, B. Zhou, C. dos Santos, Ç. glar Gulçehre, and B. Xiang, “Abstractive text summarization using sequence-to-sequence RNNs and beyond,” *CoNLL 2016*, p. 280, 2016.
- [15] A. M. Rush, S. Chopra, and J. Weston, “A neural attention model for abstractive sentence summarization,” in *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, 2015, pp. 379–389.
- [16] C.-Y. Lin, “ROUGE: A package for automatic evaluation of summaries,” *Text Summarization Branches Out*, 2004.
- [17] R. Paulus, C. Xiong, and R. Socher, “A deep reinforced model for abstractive summarization,” *arXiv preprint arXiv:1705.04304*, 2017.
- [18] A. Celikyilmaz, A. Bosselut, X. He, and Y. Choi, “Deep communicating agents for abstractive summarization,” in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, vol. 1, 2018, pp. 1662–1675.
- [19] I. Sutskever, O. Vinyals, and Q. V. Le, “Sequence to sequence learning with neural networks,” in *Advances in neural information processing systems*, 2014, pp. 3104–3112.
- [20] K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, “Learning phrase representations using RNN encoder–decoder for statistical machine translation,” in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2014, pp. 1724–1734.
- [21] T. Luong, H. Pham, and C. D. Manning, “Effective approaches to attention-based neural machine translation,” in *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, 2015, pp. 1412–1421.
- [22] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey *et al.*, “Google’s neural machine translation system: Bridging the gap between human and machine translation,” *arXiv preprint arXiv:1609.08144*, 2016.
- [23] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” *arXiv preprint arXiv:1409.0473*, 2014.
- [24] S. Shen, Y. Cheng, Z. He, W. He, H. Wu, M. Sun, and Y. Liu, “Minimum risk training for neural machine translation,” in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, vol. 1, 2016, pp. 1683–1692.
- [25] G. Klein, Y. Kim, Y. Deng, J. Senellart, and A. Rush, “OpenNMT: Open-source toolkit for neural machine translation,” *Proceedings of ACL 2017, System Demonstrations*, pp. 67–72, 2017.
- [26] S. Chopra, M. Auli, and A. M. Rush, “Abstractive sentence summarization with attentive recurrent neural networks,” in *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2016, pp. 93–98.
- [27] S.-Q. Shen, Y.-K. Lin, C.-C. Tu, Y. Zhao, Z.-Y. Liu, M.-S. Sun *et al.*, “Recent advances on neural headline generation,” *Journal of Computer Science and Technology*, vol. 32, no. 4, pp. 768–784, 2017.
- [28] D. Bahdanau, J. Chorowski, D. Serdyuk, P. Brakel, and Y. Bengio, “End-to-end attention-based large vocabulary speech recognition,” in *Acoustics, Speech and Signal Processing (ICASSP), 2016 IEEE International Conference on*. IEEE, 2016, pp. 4945–4949.
- [29] A. Graves and N. Jaitly, “Towards end-to-end speech recognition with recurrent neural networks,” in *International Conference on Machine Learning*, 2014, pp. 1764–1772.
- [30] Y. Miao, M. Gowayyed, and F. Metze, “EESSEN: End-to-end speech recognition using deep rnn models and wfst-based decoding,” in *Automatic Speech Recognition and Understanding (ASRU), 2015 IEEE Workshop on*. IEEE, 2015, pp. 167–174.
- [31] J. L. Elman, “Finding structure in time,” *Cognitive science*, vol. 14, no. 2, pp. 179–211, 1990.
- [32] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [33] K. M. Hermann, T. Kocisky, E. Grefenstette, L. Espeholt, W. Kay, M. Suleyman, and P. Blunsom, “Teaching machines to read and comprehend,” in *Advances in Neural Information Processing Systems*, 2015, pp. 1693–1701.
- [34] D. Chen, J. Bolton, and C. D. Manning, “A thorough examination of the CNN/Daily Mail reading comprehension task,” in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, vol. 1, 2016, pp. 2358–2367.
- [35] O. Vinyals, M. Fortunato, and N. Jaitly, “Pointer networks,” in *Advances in Neural Information Processing Systems*, 2015, pp. 2692–2700.
- [36] J. Gu, Z. Lu, H. Li, and V. O. Li, “Incorporating copying mechanism in sequence-to-sequence learning,”

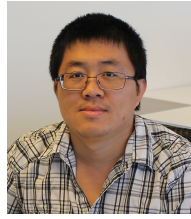
- in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, vol. 1, 2016, pp. 1631–1640.
- [37] C. Gulcehre, S. Ahn, R. Nallapati, B. Zhou, and Y. Bengio, “Pointing the unknown words,” in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, vol. 1, 2016, pp. 140–149.
- [38] W. Zeng, W. Luo, S. Fidler, and R. Urtasun, “Efficient summarization with read-again and copy mechanism,” *arXiv preprint arXiv:1611.03382*, 2016.
- [39] J. Cheng and M. Lapata, “Neural summarization by extracting sentences and words,” in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, vol. 1, 2016, pp. 484–494.
- [40] Y. Miao and P. Blunsom, “Language as a latent variable: Discrete generative models for sentence compression,” in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 2016, pp. 319–328.
- [41] K. Song, L. Zhao, and F. Liu, “Structure-infused copy mechanisms for abstractive summarization,” in *Proceedings of the 27th International Conference on Computational Linguistics*. Association for Computational Linguistics, 2018, pp. 1717–1729.
- [42] Y. Jiang and M. Bansal, “Closed-book training to improve summarization encoder memory,” in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 4067–4077.
- [43] S. Gehrmann, Y. Deng, and A. Rush, “Bottom-up abstractive summarization,” in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 4098–4109.
- [44] S. Bengio, O. Vinyals, N. Jaitly, and N. Shazeer, “Scheduled sampling for sequence prediction with recurrent neural networks,” in *Advances in Neural Information Processing Systems*, 2015, pp. 1171–1179.
- [45] A. Venkatraman, M. Hebert, and J. A. Bagnell, “Improving multi-step prediction of learned time series models,” in *Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.
- [46] M. Ranzato, S. Chopra, M. Auli, and W. Zaremba, “Sequence level training with recurrent neural networks,” *arXiv preprint arXiv:1511.06732*, 2015.
- [47] Y. Keneshloo, T. Shi, N. Ramakrishnan, and C. K. Reddy, “Deep reinforcement learning on sequence to sequence models,” *arXiv preprint arXiv:*, 2018.
- [48] Y. Bengio, R. Ducharme, P. Vincent, and C. Jauvin, “A neural probabilistic language model,” *Journal of machine learning research*, vol. 3, no. Feb, pp. 1137–1155, 2003.
- [49] R. J. Williams and D. Zipser, “A learning algorithm for continually running fully recurrent neural networks,” *Neural computation*, vol. 1, no. 2, pp. 270–280, 1989.
- [50] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, “Bleu: a method for automatic evaluation of machine translation,” in *Proceedings of the 40th annual meeting of the association for computational linguistics*. Association for Computational Linguistics, 2002, pp. 311–318.
- [51] R. J. Williams, “Simple statistical gradient-following algorithms for connectionist reinforcement learning,” in *Reinforcement Learning*. Springer, 1992, pp. 5–32.
- [52] D. Bahdanau, P. Brakel, K. Xu, A. Goyal, R. Lowe, J. Pineau, A. Courville, and Y. Bengio, “An actor-critic algorithm for sequence prediction,” *arXiv preprint arXiv:1607.07086*, 2016.
- [53] S. J. Rennie, E. Marcheret, Y. Mroueh, J. Ross, and V. Goel, “Self-critical sequence training for image captioning,” in *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*. IEEE, 2017, pp. 1179–1195.
- [54] J. Li, W. Monroe, A. Ritter, D. Jurafsky, M. Galley, and J. Gao, “Deep reinforcement learning for dialogue generation,” in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 2016, pp. 1192–1202.
- [55] J. Ling and A. Rush, “Coarse-to-fine attention models for document summarization,” in *Proceedings of the Workshop on New Frontiers in Summarization*, 2017, pp. 33–42.
- [56] X. Zhang and M. Lapata, “Sentence simplification with deep reinforcement learning,” in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 2017, pp. 584–594.
- [57] R. Pasunuru and M. Bansal, “Multi-reward reinforced summarization with saliency and entailment,” in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, vol. 2, 2018, pp. 646–653.
- [58] P. Li, L. Bing, and W. Lam, “Actor-critic based training framework for abstractive summarization,” *arXiv preprint arXiv:1803.11070*, 2018.
- [59] Y.-C. Chen and M. Bansal, “Fast abstractive summarization with reinforce-selected sentence rewriting,” in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 2018, pp. 675–686.
- [60] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, “Empirical evaluation of gated recurrent neural networks on sequence modeling,” in *NIPS 2014 Workshop on Deep Learning, December 2014*, 2014.
- [61] Y. Bengio, P. Simard, and P. Frasconi, “Learning long-term dependencies with gradient descent is difficult,” *IEEE transactions on neural networks*, vol. 5, no. 2, pp. 157–166, 1994.
- [62] R. Pascanu, T. Mikolov, and Y. Bengio, “On the difficulty of training recurrent neural networks,” in *International Conference on Machine Learning*, 2013, pp. 1310–1318.
- [63] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems*, 2017, pp. 6000–6010.

- [64] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [65] N. Kalchbrenner, L. Espeholt, K. Simonyan, A. v. d. Oord, A. Graves, and K. Kavukcuoglu, “Neural machine translation in linear time,” *arXiv preprint arXiv:1610.10099*, 2016.
- [66] J. Bradbury, S. Merity, C. Xiong, and R. Socher, “Quasi-recurrent neural networks,” *arXiv preprint arXiv:1611.01576*, 2016.
- [67] J. Gehring, M. Auli, D. Grangier, D. Yarats, and Y. N. Dauphin, “Convolutional sequence to sequence learning,” in *International Conference on Machine Learning*, 2017, pp. 1243–1252.
- [68] S. Hochreiter, Y. Bengio, P. Frasconi, J. Schmidhuber *et al.*, “Gradient flow in recurrent nets: the difficulty of learning long-term dependencies,” 2001.
- [69] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, “Wavenet: A generative model for raw audio,” in *9th ISCA Speech Synthesis Workshop*, pp. 125–125.
- [70] D. Balduzzi and M. Ghifary, “Strongly-typed recurrent neural networks,” in *Proceedings of the 33rd International Conference on International Conference on Machine Learning-Volume 48*. JMLR. org, 2016, pp. 1292–1300.
- [71] J. Gehring, M. Auli, D. Grangier, and Y. Dauphin, “A convolutional encoder model for neural machine translation,” in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, vol. 1, 2017, pp. 123–135.
- [72] A. Fan, D. Grangier, and M. Auli, “Controllable abstractive summarization,” *arXiv preprint arXiv:1711.05217*, 2017.
- [73] Y. N. Dauphin, A. Fan, M. Auli, and D. Grangier, “Language modeling with gated convolutional networks,” in *International Conference on Machine Learning*, 2017, pp. 933–941.
- [74] Q. Zhou, N. Yang, F. Wei, and M. Zhou, “Selective encoding for abstractive sentence summarization,” in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, vol. 1, 2017, pp. 1095–1104.
- [75] J. Tan, X. Wan, and J. Xiao, “Abstractive document summarization with a graph-based attentional neural model,” in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, vol. 1, 2017, pp. 1171–1181.
- [76] Y. Xia, F. Tian, L. Wu, J. Lin, T. Qin, N. Yu, and T.-Y. Liu, “Deliberation networks: Sequence generation beyond one-pass decoding,” in *Advances in Neural Information Processing Systems*, 2017, pp. 1782–1792.
- [77] P. Li, W. Lam, L. Bing, and Z. Wang, “Deep recurrent generative decoder for abstractive text summarization,” in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 2017, pp. 2091–2100.
- [78] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” *arXiv preprint arXiv:1312.6114*, 2013.
- [79] D. J. Rezende, S. Mohamed, and D. Wierstra, “Stochastic backpropagation and approximate inference in deep generative models,” in *International Conference on Machine Learning*, 2014, pp. 1278–1286.
- [80] W.-T. Hsu, C.-K. Lin, M.-Y. Lee, K. Min, J. Tang, and M. Sun, “A unified model for extractive and abstractive summarization using inconsistency loss,” *arXiv preprint arXiv:1805.06266*, 2018.
- [81] R. Nallapati, F. Zhai, and B. Zhou, “Summarunner: A recurrent neural network based sequence model for extractive summarization of documents.” 2017.
- [82] C. Li, W. Xu, S. Li, and S. Gao, “Guiding generation for abstractive text summarization based on key information guide network,” in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, vol. 2, 2018, pp. 55–60.
- [83] R. Mihalcea and P. Tarau, “Texttrank: Bringing order into text,” in *Proceedings of the 2004 conference on empirical methods in natural language processing*, 2004.
- [84] A. Cohan, F. Dernoncourt, D. S. Kim, T. Bui, S. Kim, W. Chang, and N. Goharian, “A discourse-aware attention model for abstractive summarization of long documents,” in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, vol. 2, 2018, pp. 615–621.
- [85] L. Page, S. Brin, R. Motwani, and T. Winograd, “The pagerank citation ranking: Bringing order to the web.” Stanford InfoLab, Tech. Rep., 1999.
- [86] R. Pasunuru, H. Guo, and M. Bansal, “Towards improving abstractive summarization via entailment generation,” in *Proceedings of the Workshop on New Frontiers in Summarization*, 2017, pp. 27–32.
- [87] H. Guo, R. Pasunuru, and M. Bansal, “Soft layer-specific multi-task summarization with entailment and question generation,” in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 2018, pp. 687–697.
- [88] B. McCann, N. S. Keskar, C. Xiong, and R. Socher, “The natural language Decathlon: Multitask learning as question answering,” *arXiv preprint arXiv:1806.08730*, 2018.
- [89] K. Gimpel, D. Batra, C. Dyer, and G. Shakhnarovich, “A systematic exploration of diversity in machine translation,” in *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, 2013, pp. 1100–1111.
- [90] J. Li, M. Galley, C. Brockett, J. Gao, and B. Dolan, “A diversity-promoting objective function for neural conversation models,” in *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language*



- Technologies*, 2016, pp. 110–119.
- [91] L. Bahl, P. Brown, P. De Souza, and R. Mercer, “Maximum mutual information estimation of hidden markov model parameters for speech recognition,” in *Acoustics, Speech, and Signal Processing, IEEE International Conference on ICASSP’86.*, vol. 11. IEEE, 1986, pp. 49–52.
  - [92] J. Li and D. Jurafsky, “Mutual information and diverse decoding improve neural machine translation,” *arXiv preprint arXiv:1601.00372*, 2016.
  - [93] J. Li, W. Monroe, and D. Jurafsky, “A simple, fast diverse decoding algorithm for neural generation,” *arXiv preprint arXiv:1611.08562*, 2016.
  - [94] A. K. Vijayakumar, M. Cogswell, R. R. Selvaraju, Q. Sun, S. Lee, D. Crandall, and D. Batra, “Diverse beam search: Decoding diverse solutions from neural sequence models,” *arXiv preprint arXiv:1610.02424*, 2016.
  - [95] A. Cibils, C. Musat, A. Hossman, and M. Baeriswyl, “Diverse beam search for increased novelty in abstractive summarization,” *arXiv preprint arXiv:1802.01457*, 2018.
  - [96] S. Guo and S. Sanner, “Probabilistic latent maximal marginal relevance,” in *Proceedings of the 33rd international ACM SIGIR conference on Research and development in information retrieval*. ACM, 2010, pp. 833–834.
  - [97] N. Moratanch and S. Chitrakala, “A survey on abstractive text summarization,” in *Circuit, Power and Computing Technologies (ICCPCT), 2016 International Conference on*. IEEE, 2016, pp. 1–7.
  - [98] P. K. Rachabathuni, “A survey on abstractive summarization techniques,” in *Inventive Computing and Informatics (ICICI), International Conference on*. IEEE, 2017, pp. 762–765.
  - [99] V. Dalal and L. G. Malik, “A survey of extractive and abstractive text summarization techniques,” in *Emerging Trends in Engineering and Technology (ICETET), 2013 6th International Conference on*. IEEE, 2013, pp. 109–110.
  - [100] M. Grusky, M. Naaman, and Y. Artzi, “Newsroom: A dataset of 1.3 million summaries with diverse extractive strategies,” in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, vol. 1, 2018, pp. 708–719.
  - [101] P. J. Werbos, “Backpropagation through time: what it does and how to do it,” *Proceedings of the IEEE*, vol. 78, no. 10, pp. 1550–1560, 1990.
  - [102] I. Sutskever, *Training recurrent neural networks*. University of Toronto Toronto, Ontario, Canada, 2013.
  - [103] D. Hu, “An introductory survey on attention mechanisms in nlp problems,” *arXiv preprint arXiv:1811.05544*, 2018.
  - [104] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, “Show, attend and tell: Neural image caption generation with visual attention,” in *International conference on machine learning*, 2015, pp. 2048–2057.
  - [105] Z. Tu, Z. Lu, Y. Liu, X. Liu, and H. Li, “Modeling coverage for neural machine translation,” in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, vol. 1, 2016, pp. 76–85.
  - [106] B. Sankaran, H. Mi, Y. Al-Onaizan, and A. Ittycheriah, “Temporal attention model for neural machine translation,” *CoRR*, vol. abs/1608.02927, 2016.
  - [107] Q. Chen, X. Zhu, Z. Ling, S. Wei, and H. Jiang, “Distraction-based neural networks for modeling documents,” in *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*. AAAI Press, 2016, pp. 2754–2760.
  - [108] H. Inan, K. Khosravi, and R. Socher, “Tying word vectors and word classifiers: A loss framework for language modeling,” *arXiv preprint arXiv:1611.01462*, 2016.
  - [109] J. Chung, K. Kastner, L. Dinh, K. Goel, A. C. Courville, and Y. Bengio, “A recurrent latent variable model for sequential data,” in *Advances in neural information processing systems*, 2015, pp. 2980–2988.
  - [110] C. Doersch, “Tutorial on variational autoencoders,” *arXiv preprint arXiv:1606.05908*, 2016.
  - [111] H. Bahuleyan, L. Mou, O. Vechtomova, and P. Poupard, “Variational attention for sequence-to-sequence models,” in *COLING*, 2018.
  - [112] Z. Yang, D. Yang, C. Dyer, X. He, A. Smola, and E. Hovy, “Hierarchical attention networks for document classification,” in *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2016, pp. 1480–1489.
  - [113] T. H. Haveliwala, “Topic-sensitive pagerank,” in *Proceedings of the 11th international conference on World Wide Web*. ACM, 2002, pp. 517–526.
  - [114] K. Lopyrev, “Generating news headlines with recurrent neural networks,” *arXiv preprint arXiv:1512.01712*, 2015.
  - [115] S. Takase, J. Suzuki, N. Okazaki, T. Hirao, and M. Nagata, “Neural headline generation on abstract meaning representation,” in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 2016, pp. 1054–1059.
  - [116] P. Nema, M. M. Khapra, A. Laha, and B. Ravindran, “Diversity driven attention model for query-based abstractive summarization,” in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, vol. 1, 2017, pp. 1063–1072.
  - [117] L. Liu, Y. Lu, M. Yang, Q. Qu, J. Zhu, and H. Li, “Generative adversarial network for abstractive text summarization,” *arXiv preprint arXiv:1711.09357*, 2017.
  - [118] J. Lin, X. SUN, S. Ma, and Q. Su, “Global encoding for abstractive summarization,” in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Association for Computational Linguistics, 2018, pp. 163–169.

- [119] L. Wang, J. Yao, Y. Tao, L. Zhong, W. Liu, and Q. Du, "A reinforced topic-aware convolutional sequence-to-sequence model for abstractive text summarization," *arXiv preprint arXiv:1805.03616*, 2018.
- [120] W. Kryściński, R. Paulus, C. Xiong, and R. Socher, "Improving abstraction in text summarization," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 1808–1817.
- [121] Y. Zhang, D. Y. Ding, T. Qian, C. D. Manning, and C. P. Langlotz, "Learning to summarize radiology findings," *EMNLP 2018*, p. 204, 2018.
- [122] T. L. Chung, B. Xu, Y. Liu, and C. Ouyang, "Main point generator: Summarizing with a focus," in *International Conference on Database Systems for Advanced Applications*. Springer, 2018, pp. 924–932.
- [123] X. Chen, S. Gao, C. Tao, Y. Song, D. Zhao, and R. Yan, "Iterative document representation learning towards summarization with polishing," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 4088–4097.
- [124] A. L. Berger, V. J. D. Pietra, and S. A. D. Pietra, "A maximum entropy approach to natural language processing," *Computational linguistics*, vol. 22, no. 1, pp. 39–71, 1996.
- [125] L. Weaver and N. Tao, "The optimal reward baseline for gradient-based reinforcement learning," in *Proceedings of the Seventeenth conference on Uncertainty in artificial intelligence*. Morgan Kaufmann Publishers Inc., 2001, pp. 538–545.
- [126] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 1998.
- [127] W. Zaremba and I. Sutskever, "Reinforcement learning neural turing machines-revised," *arXiv preprint arXiv:1505.00521*, 2015.
- [128] F. J. Och, "Minimum error rate training in statistical machine translation," in *Proceedings of the 41st Annual Meeting on Association for Computational Linguistics-Volume 1*. Association for Computational Linguistics, 2003, pp. 160–167.
- [129] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [130] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [131] X. Zhang, M. Lapata, F. Wei, and M. Zhou, "Neural latent extractive document summarization," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 779–784.



**Tian Shi** is working toward the Ph.D. degree in the Department of Computer Science, Virginia Tech. His research interests include deep learning, matrix factorization, text summarization, and topic modeling. He received a Ph.D. degree in Physical Chemistry from Wayne State University in 2016.



**Yaser Keneshloo** received his Masters degree in Computer Engineering from Iran University of Science and Technology in 2012. Currently, he is pursuing his Ph.D in the Department of Computer Science at Virginia Tech. His research interests includes machine learning, data mining, and deep learning.



**Naren Ramakrishnan** is the Thomas L. Phillips Professor of Engineering at Virginia Tech. He directs the Discovery Analytics Center, a university-wide effort that brings together researchers from computer science, statistics, mathematics, and electrical and computer engineering to tackle knowledge discovery problems in important areas of national interest. His work has been featured in the Wall Street Journal, Newsweek, Smithsonian Magazine, PBS/NoVA Next, Chronicle of Higher Education, and Popular Science, among other venues. He serves on the editorial boards of IEEE Computer, ACM Transactions on Knowledge Discovery from Data, Data Mining and Knowledge Discovery, IEEE Transactions on Knowledge and Data Engineering, and other journals. He received his PhD in Computer Sciences from Purdue University.



**Chandan K. Reddy** is an Associate Professor in the Department of Computer Science at Virginia Tech. He received his Ph.D. from Cornell University and M.S. from Michigan State University. His primary research interests are Data Mining and Machine Learning with applications to Healthcare Analytics and Social Network Analysis. His research is funded by the National Science Foundation, the National Institutes of Health, the Department of Transportation, and the Susan G. Komen for the Cure Foundation. He has published over 110 peer-reviewed articles in leading conferences and journals. He received several awards for his research work including the Best Application Paper Award at ACM SIGKDD conference in 2010, Best Poster Award at IEEE VAST conference in 2014, Best Student Paper Award at IEEE ICDM conference in 2016, and was a finalist of the INFORMS Franz Edelman Award Competition in 2011. He is an associate editor of the ACM Transactions on Knowledge Discovery and Data Mining and PC Co-Chair of ASONAM 2018. He is a senior member of the IEEE and life member of the ACM.