

# Selective Encoding for Abstractive Sentence Summarization

Qingyu Zhou<sup>†\*</sup> Nan Yang<sup>‡</sup> Furu Wei<sup>‡</sup> Ming Zhou<sup>‡</sup>

<sup>†</sup>Harbin Institute of Technology, Harbin, China

<sup>‡</sup>Microsoft Research, Beijing, China

qyzhou@hit.edu.cn {nanya, fuwei, mingzhou}@microsoft.com

## Abstract

We propose a selective encoding model to extend the sequence-to-sequence framework for abstractive sentence summarization. It consists of a sentence encoder, a selective gate network, and an attention equipped decoder. The sentence encoder and decoder are built with recurrent neural networks. The selective gate network constructs a second level sentence representation by controlling the information flow from encoder to decoder. The second level representation is tailored for sentence summarization task, which leads to better performance. We evaluate our model on the English Gigaword, DUC 2004 and MSR abstractive sentence summarization datasets. The experimental results show that the proposed selective encoding model outperforms the state-of-the-art baseline models.

## 1 Introduction

Sentence summarization aims to shorten a given sentence and produce a brief summary of it. This is different from document level summarization task since it is hard to apply existing techniques in extractive methods, such as extracting sentence level features and ranking sentences. Early works propose using rule-based methods (Zajic et al., 2007), syntactic tree pruning methods (Knight and Marcu, 2002), statistical machine translation techniques (Banko et al., 2000) and so on for this task. We focus on abstractive sentence summarization task in this paper.

Recently, neural network models have been applied in this task. Rush et al. (2015) use auto-constructed sentence-headline pairs to train a neu-

ral network summarization model. They use a Convolutional Neural Network (CNN) encoder and feed-forward neural network language model decoder for this task. Chopra et al. (2016) extend their work by replacing the decoder with Recurrent Neural Network (RNN). Nallapati et al. (2016) follow this line and change the encoder to RNN to make it a full RNN based sequence-to-sequence model (Sutskever et al., 2014).

the sri lankan government on wednesday announced the closure of government schools with immediate effect as a military campaign against tamil separatists escalated in the north of the country .



sri lanka closes schools as war escalates

Figure 1: An abstractive sentence summarization system may produce the output summary by distilling the salient information from the highlight to generate a fluent sentence. We model the distilling process with selective encoding.

All the above works fall into the encoding-decoding paradigm, which first encodes the input sentence to an abstract representation and then decodes the intended output sentence based on the encoded information. As an extension of the encoding-decoding framework, attention-based approach (Bahdanau et al., 2015) has been broadly used: the encoder produces a list of vectors for all tokens in the input, and the decoder uses an attention mechanism to dynamically extract encoded information and align with the output tokens. This approach achieves huge success in tasks like machine translation, where alignment between all parts of the input and output are required. However, in abstractive sentence summarization, there is no explicit alignment relationship between the input sentence and the summary ex-

\*Contribution during internship at Microsoft Research.

cept for the extracted common words. The challenge here is not to infer the alignment, but to select the highlights while filtering out secondary information in the input. A desired work-flow for abstractive sentence summarization is encoding, selection, and decoding. After selecting the important information from an encoded sentence, the decoder produces the output summary using the selected information. For example, in Figure 1, given the input sentence, the summarization system first selects the important information, and then rephrases or paraphrases to produce a well-organized summary. Although this is implicitly modeled in the encoding-decoding framework, we argue that abstractive sentence summarization shall benefit from explicitly modeling this selection process.

In this paper we propose **Selective Encoding for Abstractive Sentence Summarization (SEASS)**. We treat the sentence summarization as a three-phase task: encoding, selection, and decoding. It consists of a sentence encoder, a selective gate network, and a summary decoder. First, the sentence encoder reads the input words through an RNN unit to construct the first level sentence representation. Then the selective gate network selects the encoded information to construct the second level sentence representation. The selective mechanism controls the information flow from encoder to decoder by applying a gate network according to the sentence information, which helps improve encoding effectiveness and release the burden of the decoder. Finally, the attention-equipped decoder generates the summary using the second level sentence representation. We conduct experiments on English Gigaword, DUC 2004 and Microsoft Research Abstractive Text Compression test sets. Our SEASS model achieves 17.54 ROUGE-2 F1, 9.56 ROUGE-2 recall and 10.63 ROUGE-2 F1 on these test sets respectively, which improves performance compared to the state-of-the-art methods.

## 2 Related Work

Abstractive sentence summarization, also known as sentence compression and similar to headline generation, is used to help compress or fuse the selected sentences in extractive document summarization systems since they may inadvertently include unnecessary information. The sentence summarization task has been long connected to the headline generation task. There are some previous

methods to solve this task, such as the linguistic rule-based method (Dorr et al., 2003). As for the statistical machine learning based methods, Banko et al. (2000) apply statistical machine translation techniques by modeling headline generation as a translation task and use 8000 article-headline pairs to train the system.

Rush et al. (2015) propose leveraging news data in Annotated English Gigaword (Napoles et al., 2012) corpus to construct large scale parallel data for sentence summarization task. They propose an ABS model, which consists of an attentive Convolutional Neural Network encoder and a neural network language model (Bengio et al., 2003) decoder. On this Gigaword test set and DUC 2004 test set, the ABS model produces the state-of-the-art results. Chopra et al. (2016) extend this work, which keeps the CNN encoder but replaces the decoder with recurrent neural networks. Their experiments shows that the CNN encoder with RNN decoder model performs better than Rush et al. (2015). Nallapati et al. (2016) further change the encoder to an RNN encoder, which leads to a full RNN sequence-to-sequence model. Besides, they enrich the encoder with lexical and statistic features which play important roles in traditional feature based summarization systems, such as NER and POS tags, to improve performance. Experiments on the Gigaword and DUC 2004 test sets show that the above models achieve state-of-the-art results.

Gu et al. (2016) and Gulcehre et al. (2016) come up similar ideas that summarization task can benefit from copying words from input sentences. Gu et al. (2016) propose CopyNet to model the copying action in response generation, which also applies for summarization task. Gulcehre et al. (2016) propose a switch gate to control whether to copy from source or generate from decoder vocabulary. Zeng et al. (2016) also propose using copy mechanism and add a scalar weight on the gate of GRU/LSTM for this task. Cheng and Lapata (2016) use an RNN based encoder-decoder for extractive summarization of documents.

Yu et al. (2016) propose a segment to segment neural transduction model for sequence-to-sequence framework. The model introduces a latent segmentation which determines correspondences between tokens of the input sequence and the output sequence. Experiments on this task show that the proposed transduction model per-

forms comparable to the ABS model. Shen et al. (2016) propose to apply Minimum Risk Training (MRT) in neural machine translation to directly optimize the evaluation metrics. Ayana et al. (2016) apply MRT on abstractive sentence summarization task and the results show that optimizing for ROUGE improves the test performance.

### 3 Problem Formulation

For sentence summarization, given an input sentence  $x = (x_1, x_2, \dots, x_n)$ , where  $n$  is the sentence length,  $x_i \in \mathcal{V}_s$  and  $\mathcal{V}_s$  is the source vocabulary, the system summarizes  $x$  by producing  $y = (y_1, y_2, \dots, y_l)$ , where  $l \leq n$  is the summary length,  $y_i \in \mathcal{V}_t$  and  $\mathcal{V}_t$  is the target vocabulary.

If  $|y| \subseteq |x|$ , which means all words in summary  $y$  must appear in given input, we denote this as *extractive* sentence summarization. If  $|y| \not\subseteq |x|$ , which means not all words in summary come from input sentence, we denote this as *abstractive* sentence summarization. Table 1 provides an example. We focus on *abstractive* sentence summarization task in this paper.

<b>Input:</b>	South Korean President Kim Young-Sam left here Wednesday on a week - long state visit to Russia and Uzbekistan for talks on North Korea 's nuclear confrontation and ways to strengthen bilateral ties .
<b>Output:</b>	Kim leaves for Russia for talks on NKorea nuclear standoff

Table 1: An abstractive sentence summarization example.

## 4 Model

As shown in Figure 2, our model consists of a sentence encoder using the Gated Recurrent Unit (GRU) (Cho et al., 2014), a selective gate network and an attention-equipped GRU decoder. First, the bidirectional GRU encoder reads the input words  $x = (x_1, x_2, \dots, x_n)$  and builds its representation  $(h_1, h_2, \dots, h_n)$ . Then the selective gate selects and filters the word representations according to the sentence meaning representation to produce a tailored sentence word representation for abstractive sentence summarization task. Lastly, the GRU decoder produces the output summary with attention to the tailored representation. In the following sections, we introduce the sentence encoder, the selective mechanism, and the summary decoder respectively.

### 4.1 Sentence Encoder

The role of the sentence encoder is to read the input sentence and construct the basic sentence representation. Here we employ a bidirectional GRU (BiGRU) as the recurrent unit, where GRU is defined as:

$$z_i = \sigma(\mathbf{W}_z[x_i, h_{i-1}]) \quad (1)$$

$$r_i = \sigma(\mathbf{W}_r[x_i, h_{i-1}]) \quad (2)$$

$$\tilde{h}_i = \tanh(\mathbf{W}_h[x_i, r_i \odot h_{i-1}]) \quad (3)$$

$$h_i = (1 - z_i) \odot h_{i-1} + z_i \odot \tilde{h}_i \quad (4)$$

where  $\mathbf{W}_z$ ,  $\mathbf{W}_r$  and  $\mathbf{W}_h$  are weight matrices.

The BiGRU consists of a forward GRU and a backward GRU. The forward GRU reads the input sentence word embeddings from left to right and gets a sequence of hidden states,  $(\vec{h}_1, \vec{h}_2, \dots, \vec{h}_n)$ . The backward GRU reads the input sentence embeddings reversely, from right to left, and results in another sequence of hidden states,  $(\vec{h}_1, \vec{h}_2, \dots, \vec{h}_n)$ :

$$\vec{h}_i = \text{GRU}(x_i, \vec{h}_{i-1}) \quad (5)$$

$$\vec{h}_i = \text{GRU}(x_i, \vec{h}_{i+1}) \quad (6)$$

The initial states of the BiGRU are set to zero vectors, i.e.,  $\vec{h}_1 = 0$  and  $\vec{h}_n = 0$ . After reading the sentence, the forward and backward hidden states are concatenated, i.e.,  $h_i = [\vec{h}_i; \vec{h}_i]$ , to get the basic sentence representation.

### 4.2 Selective Mechanism

In the sequence-to-sequence machine translation (MT) model, the encoder and decoder are responsible for mapping input sentence information to a list of vectors and decoding the sentence representation vectors to generate an output sentence (Bahdanau et al., 2015). Some previous works apply this framework to summarization generation tasks (Nallapati et al., 2016; Gu et al., 2016; Gulcehre et al., 2016). However, abstractive sentence summarization is different from MT in two ways. First, there is no explicit alignment relationship between the input sentence and the output summary except for the common words. Second, summarization task needs to keep the highlights and remove the unnecessary information, while MT needs to keep all information literally.

Herein, we propose a selective mechanism to model the selection process for abstractive sentence summarization. The selective mechanism

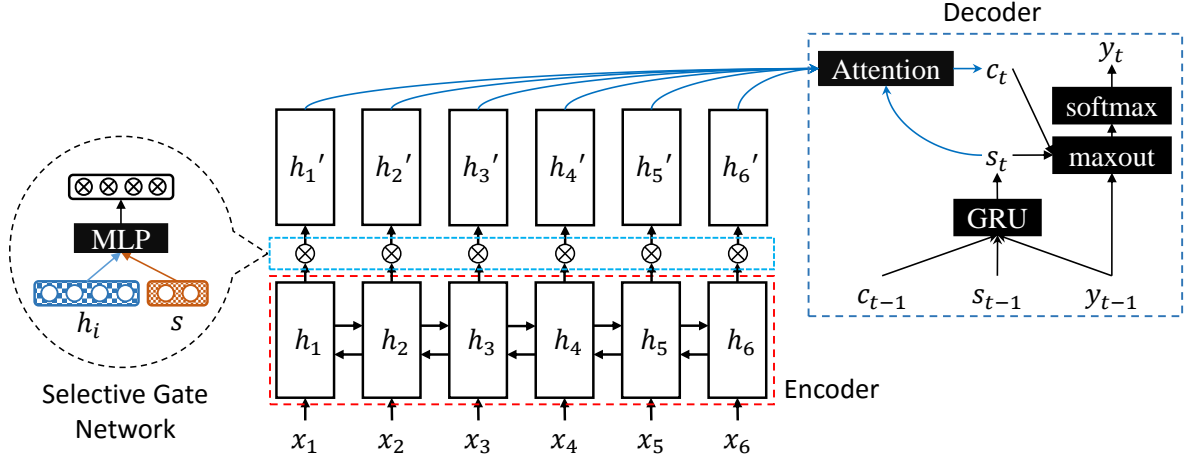


Figure 2: Overview of the **Selective Encoding for Abstractive Sentence Summarization (SEASS)**.

extends the sequence-to-sequence model by constructing a tailored representation for abstractive sentence summarization task. Concretely, the selective gate network in our model takes two vector inputs, the sentence word vector  $h_i$  and the sentence representation vector  $s$ . The sentence word vector  $h_i$  is the output of the BiGRU encoder and represents the meaning and context information of word  $x_i$ . The sentence vector  $s$  is used to represent the meaning of the sentence. For each word  $x_i$ , the selective gate network generates a gate vector  $sGate_i$  using  $h_i$  and  $s$ , then the tailored representation is constructed, i.e.,  $h'_i$ .

In detail, we concatenate the last forward hidden state  $\vec{h}_n$  and backward hidden state  $\overleftarrow{h}_1$  as the sentence representation  $s$ :

$$s = \begin{bmatrix} \vec{h}_n \\ \overleftarrow{h}_1 \end{bmatrix} \quad (7)$$

For each time step  $i$ , the selective gate takes the sentence representation  $s$  and BiGRU hidden  $h_i$  as inputs to compute the gate vector  $sGate_i$ :

$$sGate_i = \sigma(\mathbf{W}_s h_i + \mathbf{U}_s s + b) \quad (8)$$

$$h'_i = h_i \odot sGate_i \quad (9)$$

where  $\mathbf{W}_s$  and  $\mathbf{U}_s$  are weight matrices,  $b$  is the bias vector,  $\sigma$  denotes sigmoid activation function, and  $\odot$  is element-wise multiplication. After the selective gate network, we obtain another sequence of vectors  $(h'_1, h'_2, \dots, h'_n)$ . This new sequence is then used as the input sentence representation for the decoder to generate the summary.

### 4.3 Summary Decoder

On top of the sentence encoder and the selective gate network, we use GRU with attention as the decoder to produce the output summary.

At each decoding time step  $t$ , the GRU reads the previous word embedding  $w_{t-1}$  and previous context vector  $c_{t-1}$  as inputs to compute the new hidden state  $s_t$ . To initialize the GRU hidden state, we use a linear layer with the last backward encoder hidden state  $\overleftarrow{h}_1$  as input:

$$s_t = \text{GRU}(w_{t-1}, c_{t-1}, s_{t-1}) \quad (10)$$

$$s_0 = \tanh(\mathbf{W}_d \overleftarrow{h}_1 + b) \quad (11)$$

where  $\mathbf{W}_d$  is the weight matrix and  $b$  is the bias vector.

The context vector  $c_t$  for current time step  $t$  is computed through the concatenate attention mechanism (Luong et al., 2015), which matches the current decoder state  $s_t$  with each encoder hidden state  $h'_i$  to get an importance score. The importance scores are then normalized to get the current context vector by weighted sum:

$$e_{t,i} = v_a^\top \tanh(\mathbf{W}_a s_{t-1} + \mathbf{U}_a h'_i) \quad (12)$$

$$\alpha_{t,i} = \frac{\exp(e_{t,i})}{\sum_{i=1}^n \exp(e_{t,i})} \quad (13)$$

$$c_t = \sum_{i=1}^n \alpha_{t,i} h'_i \quad (14)$$

We then combine the previous word embedding  $w_{t-1}$ , the current context vector  $c_t$ , and the decoder state  $s_t$  to construct the readout state  $r_t$ . The readout state is then passed through a maxout hidden layer (Goodfellow et al., 2013) to predict the

next word with a softmax layer over the decoder vocabulary.

$$r_t = \mathbf{W}_r w_{t-1} + \mathbf{U}_r c_t + \mathbf{V}_r s_t \quad (15)$$

$$m_t = [\max\{r_{t,2j-1}, r_{t,2j}\}]_{j=1,\dots,d}^\top \quad (16)$$

$$p(y_t|y_1, \dots, y_{t-1}) = \text{softmax}(\mathbf{W}_o m_t) \quad (17)$$

where  $\mathbf{W}_a$ ,  $\mathbf{U}_a$ ,  $\mathbf{W}_r$ ,  $\mathbf{U}_r$ ,  $\mathbf{V}_r$  and  $\mathbf{W}_o$  are weight matrices. Readout state  $r_t$  is a  $2d$ -dimensional vector, and the maxout layer (Equation 16) picks the max value for every two numbers in  $r_t$  and produces a  $d$ -dimensional vector  $m_t$ .

#### 4.4 Objective Function

Our goal is to maximize the output summary probability given the input sentence. Therefore, we optimize the negative log-likelihood loss function:

$$J(\theta) = -\frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \log p(y|x) \quad (18)$$

where  $\mathcal{D}$  denotes a set of parallel sentence-summary pairs and  $\theta$  is the model parameter. We use Stochastic Gradient Descent (SGD) with mini-batch to learn the model parameter  $\theta$ .

## 5 Experiments

In this section we introduce the dataset we use, the evaluation metric, the implementation details, the baselines we compare to, and the performance of our system.

### 5.1 Dataset

**Training Set** For our training set, we use a parallel corpus which is constructed from the Annotated English Gigaword dataset (Napoles et al., 2012) as mentioned in Rush et al. (2015). The parallel corpus is produced by pairing the first sentence and the headline in the news article with some heuristic rules. We use the script<sup>1</sup> released by Rush et al. (2015) to pre-process and extract the training and development datasets. The script performs various basic text normalization, including PTB tokenization, lower-casing, replacing all digit characters with #, and replacing word types seen less than 5 times with  $\langle unk \rangle$ . The extracted corpus contains about 3.8M sentence-summary pairs for the training set and 189K examples for the development set.

For our test set, we use the English Gigaword, DUC 2004, and Microsoft Research Abstractive Text Compression test sets.

<sup>1</sup><https://github.com/facebook/NAMAS>

**English Gigaword Test Set** We randomly sample 8000 pairs from the extracted development set as our development set since it is relatively large. For the test set, we use the same randomly held-out test set of 2000 sentence-summary pairs as Rush et al. (2015).<sup>2</sup>

We also find that except for the empty titles, this test set has some invalid lines like the input sentence containing only one word. Therefore, we further sample 2000 pairs as our internal test set and release it for future works<sup>3</sup>.

**DUC 2004 Test Set** We employ DUC 2004 data for tasks 1 & 2 (Over et al., 2007) in our experiments as one of the test sets since it is too small to train a neural network model on. The dataset pairs each document with 4 different human-written reference summaries which are capped at 75 bytes. It has 500 input sentences with each sentence paired with 4 summaries.

**MSR-ATC Test Set** Toutanova et al. (2016) release a new dataset for sentence summarization task by crowdsourcing. This dataset contains approximately 6,000 source text sentences with multiple manually-created summaries (about 26,000 sentence-summary pairs in total). Toutanova et al. (2016) provide a standard split of the data into training, development, and test sets, with 4,936, 448 and 785 input sentences respectively. Since the training set is too small, we only use the test set as one of our test sets. We denote this dataset as MSR-ATC (Microsoft Research Abstractive Text Compression) test set in the following.

Table 2 summarizes the statistic information of the three datasets we used.

### 5.2 Evaluation Metric

We employ ROUGE (Lin, 2004) as our evaluation metric. ROUGE measures the quality of summary by computing overlapping lexical units, such as unigram, bigram, trigram, and longest common subsequence (LCS). It becomes the standard evaluation metric for DUC shared tasks and popular for summarization evaluation. Following previous work, we use ROUGE-1 (unigram), ROUGE-2 (bi-

<sup>2</sup>Thanks to Rush et al. (2015), we acquired the test set they used. Following Chopra et al. (2016), we remove pairs with empty titles resulting in slightly different accuracy compared to Rush et al. (2015) for their systems. The cleaned test set contains 1951 sentence-summary pairs.

<sup>3</sup>Our development and test sets can be found at <https://res.qyzhou.me>

Data Set	Giga	DUC <sup>†</sup>	MSR <sup>†</sup>
#(sent)	3.99M	500	785
#(sentWord)	125M	17.8K	29K
#(summWord)	33M	20.9K	85.9K
#(ref)	1	4	3-5
AvgInputLen	31.35	35.56	36.97
AvgSummLen	8.23	10.43	25.5

Table 2: Data statistics for the English Gigaword, DUC 2004 and MSR-ATC datasets.  $\#(x)$  denotes the number of  $x$ , e.g.,  $\#(\text{ref})$  is the number of reference summaries of an input sentence. AvgInputLen is the average input sentence length and AvgSummLen is the average summary length. <sup>†</sup>DUC 2004 and MSR-ATC datasets are for test purpose only.

gram) and ROUGE-L (LCS) as the evaluation metrics in the reported experimental results.

### 5.3 Implementation Details

**Model Parameters** The input and output vocabularies are collected from the training data, which have 119,504 and 68,883 word types respectively. We set the word embedding size to 300 and all GRU hidden state sizes to 512. We use dropout (Srivastava et al., 2014) with probability  $p = 0.5$ .

**Model Training** We initialize model parameters randomly using a Gaussian distribution with Xavier scheme (Glorot and Bengio, 2010). We use Adam (Kingma and Ba, 2015) as our optimizing algorithm. For the hyperparameters of Adam optimizer, we set the learning rate  $\alpha = 0.001$ , two momentum parameters  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$  respectively, and  $\epsilon = 10^{-8}$ . During training, we test the model performance (ROUGE-2 F1) on development set for every 2,000 batches. We halve the Adam learning rate  $\alpha$  if the ROUGE-2 F1 score drops for twelve consecutive tests on development set. We also apply gradient clipping (Pascanu et al., 2013) with range  $[-5, 5]$  during training. To both speed up the training and converge quickly, we use mini-batch size 64 by grid search.

**Beam Search** We use beam search to generate multiple summary candidates to get better results. To avoid favoring shorter outputs, we average the ranking score along the beam path by dividing it by the number of generated words. To both decode fast and get better results, we set the beam size to

12 in our experiments.

### 5.4 Baseline

We compare SEASS model with the following state-of-the-art baselines:

**ABS** Rush et al. (2015) use an attentive CNN encoder and NNLM decoder to do the sentence summarization task. We trained this baseline model with the released code<sup>1</sup> and evaluate it with our internal English Gigaword test set and MSR-ATC test set.

**ABS+** Based on ABS model, Rush et al. (2015) further tune their model using DUC 2003 dataset, which leads to improvements on DUC 2004 test set.

**CAs2s** As an extension of the ABS model, Chopra et al. (2016) use a convolutional attention-based encoder and RNN decoder, which outperforms the ABS model.

**Feats2s** Nallapati et al. (2016) use a full RNN sequence-to-sequence encoder-decoder model and add some features to enhance the encoder, such as POS tag, NER, and so on.

**Luong-NMT** Neural machine translation model of Luong et al. (2015) with two-layer LSTMs for the encoder-decoder with 500 hidden units in each layer implemented in (Chopra et al., 2016).

**s2s+att** We also implement a sequence-to-sequence model with attention as our baseline and denote it as “s2s+att”.

### 5.5 Results

We report ROUGE F1, ROUGE recall and ROUGE F1 for English Gigaword, DUC 2004 and MSR-ATC test sets respectively. We use the official ROUGE script (version 1.5.5)<sup>4</sup> to evaluate the summarization quality in our experiments. For English Gigaword<sup>5</sup> and MSR-ATC<sup>6</sup> test sets, the outputs have different lengths so we evaluate the system with F1 metric. As for the DUC 2004 test set<sup>7</sup>, the task requires the system to produce a fixed length summary (75 bytes), therefore we employ ROUGE recall as the evaluation metric. To satisfy the length requirement, we decode the output summary to a roughly expected length following Rush et al. (2015).

<sup>4</sup><http://www.berouge.com/>

<sup>5</sup>The ROUGE evaluation option is the same as Rush et al. (2015), -m -n 2 -w 1.2

<sup>6</sup>The ROUGE evaluation option is, -m -n 2 -w 1.2

<sup>7</sup>The ROUGE evaluation option is, -m -b 75 -n 2 -w 1.2

**English Gigaword** We acquire the test set from [Rush et al. \(2015\)](#) so we can make fair comparisons to the baselines.

Models	RG-1	RG-2	RG-L
ABS (beam) <sup>‡</sup>	29.55 <sup>·</sup>	11.32 <sup>·</sup>	26.42 <sup>·</sup>
ABS+ (beam) <sup>‡</sup>	29.76 <sup>·</sup>	11.88 <sup>·</sup>	26.96 <sup>·</sup>
Feats2s (beam) <sup>‡</sup>	32.67 <sup>·</sup>	15.59 <sup>·</sup>	30.64 <sup>·</sup>
CAs2s (greedy) <sup>‡</sup>	33.10 <sup>·</sup>	14.45 <sup>·</sup>	30.25 <sup>·</sup>
CAs2s (beam) <sup>‡</sup>	33.78 <sup>·</sup>	15.97 <sup>·</sup>	31.15 <sup>·</sup>
Luong-NMT (beam) <sup>‡</sup>	33.10 <sup>·</sup>	14.45 <sup>·</sup>	30.71 <sup>·</sup>
s2s+att (greedy)	33.18 <sup>·</sup>	14.79 <sup>·</sup>	30.80 <sup>·</sup>
s2s+att (beam)	34.04 <sup>·</sup>	15.95 <sup>·</sup>	31.68 <sup>·</sup>
SEASS (greedy)	35.48	16.50	32.93
SEASS (beam)	<b>36.15</b>	<b>17.54</b>	<b>33.63</b>

Table 3: Full length ROUGE F1 evaluation results on the English Gigaword test set used by [Rush et al. \(2015\)](#). RG in the Table denotes ROUGE. Results with <sup>‡</sup> mark are taken from the corresponding papers. The superscript <sup>·</sup> indicates that our SEASS model with beam search performs significantly better than it as given by the 95% confidence interval in the official ROUGE script.

Models	RG-1	RG-2	RG-L
ABS (beam)	37.41 <sup>·</sup>	15.87 <sup>·</sup>	34.70 <sup>·</sup>
s2s+att (greedy)	42.41 <sup>·</sup>	20.76 <sup>·</sup>	39.84 <sup>·</sup>
s2s+att (beam)	43.76 <sup>·</sup>	22.28 <sup>·</sup>	41.14 <sup>·</sup>
SEASS (greedy)	45.27	22.88	42.20
SEASS (beam)	<b>46.86</b>	<b>24.58</b>	<b>43.53</b>

Table 4: Full length ROUGE F1 evaluation on our internal English Gigaword test data. The superscript <sup>·</sup> indicates that our SEASS model performs significantly better than it as given by the 95% confidence interval in the official ROUGE script.

In Table 3, we report the ROUGE F1 score of our model and the baseline methods. Our SEASS model with beam search outperforms all baseline models by a large margin. Even for greedy search, our model still performs better than other methods which used beam search. For the popular ROUGE-2 metric, our SEASS model achieves 17.54 F1 score and performs better than the previous works. Compared to the ABS model, our model has a 6.22 ROUGE-2 F1 relative gain. Compared to the highest CAs2s baseline, our model achieves 1.57

ROUGE-2 F1 improvement and passes the significant test according to the official ROUGE script.

Table 4 summarizes our results on our internal test set using ROUGE F1 evaluation metrics. The performance on our internal test set is comparable to our development set, which achieves 24.58 ROUGE-2 F1 and outperforms the baselines.

**DUC 2004** We evaluate our model using the ROUGE recall score since the reference summaries of the DUC 2004 test set are capped at 75 bytes. Therefore, we decode the summary to a fixed length 18 to ensure that the generated summary satisfies the minimum length requirement. As summarized in Table 5, our SEASS outperforms all the baseline methods and achieves 29.21, 9.56 and 25.51 for ROUGE 1, 2 and L recall. Compared to the ABS+ model which is tuned using DUC 2003 data, our model performs significantly better by 1.07 ROUGE-2 recall score and is trained only with English Gigaword sentence-summary data without being tuned using DUC data.

Models	RG-1	RG-2	RG-L
ABS (beam) <sup>‡</sup>	26.55 <sup>·</sup>	7.06 <sup>·</sup>	22.05 <sup>·</sup>
ABS+ (beam) <sup>‡</sup>	28.18 <sup>·</sup>	8.49 <sup>·</sup>	23.81 <sup>·</sup>
Feats2s (beam) <sup>‡</sup>	28.35 <sup>·</sup>	9.46	24.59 <sup>·</sup>
CAs2s (greedy) <sup>‡</sup>	29.13	7.62 <sup>·</sup>	23.92 <sup>·</sup>
CAs2s (beam) <sup>‡</sup>	28.97	8.26 <sup>·</sup>	24.06 <sup>·</sup>
Luong-NMT (beam) <sup>‡</sup>	28.55	8.79 <sup>·</sup>	24.43 <sup>·</sup>
s2s+att (greedy)	27.03 <sup>·</sup>	7.89 <sup>·</sup>	23.80 <sup>·</sup>
s2s+att (beam)	28.13	9.25	24.76
SEASS (greedy)	28.68	8.55	25.04
SEASS (beam)	<b>29.21</b>	<b>9.56</b>	<b>25.51</b>

Table 5: ROUGE recall evaluation results on DUC 2004 test set. All these models are tested using beam search. Results with <sup>‡</sup> mark are taken from the corresponding papers. The superscript <sup>·</sup> indicates that our SEASS model performs significantly better than it as given by the 95% confidence interval in the official ROUGE script.

**MSR-ATC** We report the full length ROUGE F1 score on the MSR-ATC test set in Table 6. To the best of our knowledge, this is the first work that reports ROUGE metric scores on the MSR-ATC dataset. Note that we only compare our model with ABS since the others are not publicly available. Our SEASS achieves 10.63 ROUGE-2 F1 and outperforms the s2s+att baseline by 1.02 points.

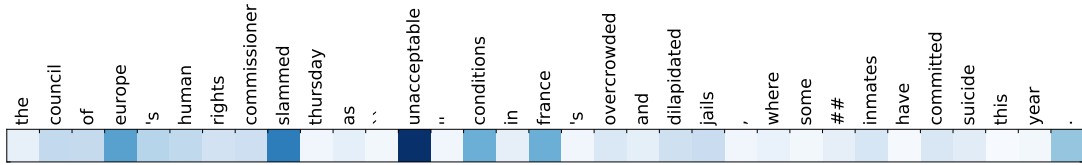


Figure 3: First derivative heat map of the output with respect to the selective gate. The important words are selected in the input sentence, such as “europe”, “slammed” and “unacceptable”. The output summary of our system is “council of europe slams french prison conditions” and the true summary is “council of europe again slams french prison conditions”.

Models	RG-1	RG-2	RG-L
ABS (beam)	20.27 <sup>*</sup>	5.26 <sup>*</sup>	17.10 <sup>*</sup>
s2s+att (greedy)	15.15 <sup>*</sup>	4.48 <sup>*</sup>	13.62 <sup>*</sup>
s2s+att (beam)	22.65 <sup>*</sup>	9.61 <sup>*</sup>	21.39 <sup>*</sup>
SEASS (greedy)	19.77	6.44	17.36
SEASS (beam)	<b>25.75</b>	<b>10.63</b>	<b>22.90</b>

Table 6: Full length ROUGE F1 evaluation on MSR-ATC test set. Beam search are used in both the baselines and our method. The superscript <sup>\*</sup> indicates that our SEASS model performs significantly better than it as given by the 95% confidence interval in the official ROUGE script.

## 6 Discussion

In this section, we first compare the performance of SEASS with the s2s+att baseline model to illustrate that the proposed method succeeds in selecting information and building tailored representation for abstractive sentence summarization. We then analyze selective encoding by visualizing the heat map.

**Effectiveness of Selective Encoding** We further test the SEASS model with different sentence lengths on English Gigaword test sets, which are merged from the Rush et al. (2015) test set and our internal test set. The length of sentences in the test sets ranges from 10 to 80. We group the sentences with an interval of 4 and get 18 different groups and we draw the first 14 groups. We find that the performance curve of our SEASS model always appears to be on the top of that of s2s+att with a certain margin. For the groups of 16, 20, 24, 32, 56 and 60, the SEASS model obtains big improvements compared to the s2s+att model. Overall, these improvements on all groups indicate that the selective encoding method benefits the abstractive sentence summarization task.

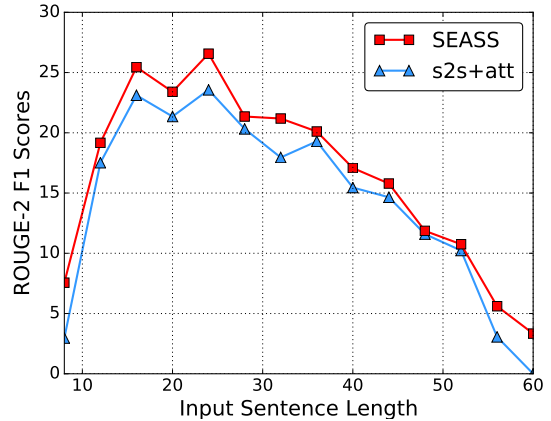


Figure 4: ROUGE-2 F1 score on different groups of input sentences in terms of their length for s2s+att baseline and our SEASS model on English Gigaword test sets.

**Saliency Heat Map of Selective Gate** Since the output of the selective gate network is a high dimensional vector, it is hard to visualize all the gate values. We use the method in Li et al. (2016) to visualize the contribution of the selective gate to the final output, which can be approximated by the first derivative. Given sentence words  $x$  with associated output summary  $y$ , the trained model associates the pair  $(x, y)$  with a score  $S_y(x)$ . The goal is to decide which gate  $g$  associated with a specific word makes the most significant contribution to  $S_y(x)$ . We approximate the  $S_y(g)$  by computing the first-order Taylor expansion since the score  $S_y(x)$  is a highly non-linear function in the deep neural network models:

$$S_y(g) \approx w(g)^T g + b \quad (19)$$

where  $w(g)$  is first the derivative of  $S_y$  with respect to the gate  $g$ :

$$w(g) = \frac{\partial(S_y)}{\partial g} \Big|_g \quad (20)$$



We then draw the Euclidean norm of the first derivative of the output  $y$  with respect to the selective gate  $g$  associated with each input words.

Figure 3 shows an example of the first derivative heat map, in which most of the important words are selected by the selective gate such as “eu-rope”, “slammed”, “unacceptable”, “conditions”, and “france”. We can observe that the selective gate determines the importance of each word before decoder, which releases the burden of it by providing tailored sentence encoding.

## 7 Conclusion

This paper proposes a selective encoding model which extends the sequence-to-sequence model for abstractive sentence summarization task. The selective mechanism mimics one of the human summarizers’ behaviors, selecting important information before writing down the summary. With the proposed selective mechanism, we build an end-to-end neural network summarization model which consists of three phases: encoding, selection, and decoding. Experimental results show that the selective encoding model greatly improves the performance with respect to the state-of-the-art methods on English Gigaword, DUC 2004 and MSR-ATC test sets.

## Acknowledgments

We thank Chuanqi Tan, Junwei Bao, Shuangzhi Wu and the anonymous reviewers for their helpful comments. We also thank Alexander M. Rush for providing the dataset for comparison and helpful discussions.

## References

Ayana, Shiqi Shen, Zhiyuan Liu, and Maosong Sun. 2016. Neural headline generation with minimum risk training. *CoRR* abs/1604.01904.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *Proceedings of 3rd International Conference for Learning Representations*. San Diego.

Michele Banko, Vibhu O Mittal, and Michael J Witbrock. 2000. Headline generation based on statistical translation. In *Proceedings of the 38th Annual Meeting on Association for Computational Linguistics*. Association for Computational Linguistics, pages 318–325.

Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. 2003. A neural probabilistic language model. *journal of machine learning research* 3(Feb):1137–1155.

Jianpeng Cheng and Mirella Lapata. 2016. Neural summarization by extracting sentences and words. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Berlin, Germany, pages 484–494.

Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using rnn encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Doha, Qatar, pages 1724–1734.

Sumit Chopra, Michael Auli, and Alexander M. Rush. 2016. 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*. Association for Computational Linguistics, San Diego, California, pages 93–98.

Bonnie Dorr, David Zajic, and Richard Schwartz. 2003. Hedge trimmer: A parse-and-trim approach to headline generation. In *Proceedings of the HLT-NAACL 03 on Text summarization workshop-Volume 5*. Association for Computational Linguistics, pages 1–8.

Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Aistats*. volume 9, pages 249–256.

Ian J Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron C Courville, and Yoshua Bengio. 2013. Max-out networks. *ICML (3)* 28:1319–1327.

Jiatao Gu, Zhengdong Lu, Hang Li, and Victor O.K. Li. 2016. 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)*. Association for Computational Linguistics, Berlin, Germany, pages 1631–1640.

Caglar Gulcehre, Sungjin Ahn, Ramesh Nallapati, Bowen Zhou, and Yoshua Bengio. 2016. Pointing the unknown words. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Berlin, Germany, pages 140–149.

Diederik Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *Proceedings of 3rd International Conference for Learning Representations*. San Diego.

- Kevin Knight and Daniel Marcu. 2002. Summarization beyond sentence extraction: A probabilistic approach to sentence compression. *Artificial Intelligence* 139(1):91–107.
- Jiwei Li, Xinlei Chen, Eduard Hovy, and Dan Jurafsky. 2016. Visualizing and understanding neural models in nlp. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, San Diego, California, pages 681–691.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out: Proceedings of the ACL-04 workshop*. Barcelona, Spain, volume 8.
- Thang Luong, Hieu Pham, and Christopher D. Manning. 2015. Effective approaches to attention-based neural machine translation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Lisbon, Portugal, pages 1412–1421.
- Ramesh Nallapati, Bowen Zhou, Çağlar Gulçehre, and Bing Xiang. 2016. Abstractive text summarization using sequence-to-sequence rnns and beyond. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*.
- Courtney Napoles, Matthew Gormley, and Benjamin Van Durme. 2012. Annotated gigaword. In *Proceedings of the Joint Workshop on Automatic Knowledge Base Construction and Web-scale Knowledge Extraction*. Association for Computational Linguistics, Stroudsburg, PA, USA, AKBC-WEKEX '12, pages 95–100.
- Paul Over, Hoa Dang, and Donna Harman. 2007. Duc in context. *Information Processing & Management* 43(6):1506–1520.
- Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. 2013. On the difficulty of training recurrent neural networks. *ICML (3)* 28:1310–1318.
- Alexander M. Rush, Sumit Chopra, and Jason Weston. 2015. A neural attention model for abstractive sentence summarization. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Lisbon, Portugal, pages 379–389.
- Shiqi Shen, Yong Cheng, Zhongjun He, Wei He, Hua Wu, Maosong Sun, and Yang Liu. 2016. Minimum risk training for neural machine translation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Berlin, Germany, pages 1683–1692.
- Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research* 15(1):1929–1958.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*. pages 3104–3112.
- Kristina Toutanova, Chris Brockett, Ke M. Tran, and Saleema Amershi. 2016. A dataset and evaluation metrics for abstractive compression of sentences and short paragraphs. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Austin, Texas, pages 340–350.
- Lei Yu, Jan Buys, and Phil Blunsom. 2016. Online segment to segment neural transduction. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Austin, Texas, pages 1307–1316.
- David Zajic, Bonnie J Dorr, Jimmy Lin, and Richard Schwartz. 2007. Multi-candidate reduction: Sentence compression as a tool for document summarization tasks. *Information Processing & Management* 43(6):1549–1570.
- Wenyuan Zeng, Wenjie Luo, Sanja Fidler, and Raquel Urtasun. 2016. Efficient summarization with read-again and copy mechanism. *arXiv preprint arXiv:1611.03382*.