

# Fast Abstractive Summarization with Reinforce-Selected Sentence Rewriting

Anonymous ACL submission

## Abstract

Inspired from how humans summarize long documents, we propose an accurate and fast summarization model that first selects salient sentences and then rewrites them abstractively (i.e., compresses and paraphrases) to generate a concise overall summary. We use a novel sentence-level policy gradient method with metric rewards to bridge the non-differentiable computation between these two neural networks in a hierarchical way, while maintaining language fluency. Empirically, we achieve the new state-of-the-art on all metrics on the CNN/Daily Mail dataset. Moreover, by first operating at the sentence-level and then the word-level, we enable parallel decoding of our neural generative model that results in substantially faster (10-20x) inference speed as well as 4x faster training speed than previous long-paragraph encoder-decoder models. We also demonstrate the generalization of our model on the test-only DUC-2002 dataset, where we achieve higher scores than a state-of-the-art model.

## 1 Introduction

The task of document summarization has two main paradigms: extractive and abstractive. The former method directly chooses and outputs the salient sentences (or phrases) in the original document (Jing and McKeown, 2000; Knight and Marcu, 2000; Martins and Smith, 2009; Berg-Kirkpatrick et al., 2011). The latter abstractive approach involves rewriting the summary (Banko et al., 2000; Zajic et al., 2004), and has seen substantial recent gains due to neural sequence-to-sequence models (Chopra et al., 2016; Nallap-

ati et al., 2016; See et al., 2017; Paulus et al., 2017). Abstractive models can be more concise by performing generation from scratch, but they suffer from slow and inaccurate encoding of very long documents, with the attention model being required to look at all encoded words (in long paragraphs) for decoding each generated summary word (slow, one by one sequentially). Abstractive models also suffer from redundancy (repetitions), especially when generating multi-sentence text.

To address both these issues and combine the advantages of both paradigms, we propose a hybrid extractive-abstractive architecture, with policy-based reinforcement learning to bridge together the two networks. Similar to how humans summarize long documents, our model first uses an extractor agent to select salient sentences or highlights, and then employs an abstractor network to rewrite (i.e., compress and paraphrase) each of these extracted sentences. To overcome the non-differentiable behavior of our extractor and train on available document-summary pairs without saliency label, we next use actor-critic policy gradient with sentence-level metric rewards to connect these two neural networks and to learn sentence saliency. We also avoid common language fluency issues (Paulus et al., 2017) by preventing the policy gradients from affecting the abstractive summarizer’s word-level training, which also stabilizes the learning dynamics. Our sentence-level RL takes into account the word-sentence hierarchy, which better models the language structure and makes parallelization possible. Our extractor combines reinforcement learning and pointer networks, which is inspired from Bello et al. (2017)’s attempt to solve the Traveling Salesman Problem. Our abstractor is a simple encoder-aligner-decoder model (with copying) and is trained on pseudo document-summary sentence pairs obtained via simple auto-

100           matic matching criteria.  
 101

102           Thus, our method incorporates the abstractive  
 103           paradigm’s advantages of concisely rewriting sentences and generating novel words from the full  
 104           vocabulary, yet it adopts intermediate extractive  
 105           behavior to improve the overall model’s quality, speed, and stability. Instead of encoding and at-  
 106           tending to every word in the long input document sequentially, our model adopts a human-inspired  
 107           coarse-to-fine approach that first extracts all the  
 108           salient sentences and then decodes (rewrites) them (in parallel). This also avoids almost all redundancy issues because the model has already chosen non-redundant salient sentences to abstractively summarize (but adding an optional final reranker component does give additional gains by removing the fewer across-sentence repetitions).

117           Empirically, our approach is the new state-of-  
 118           the-art on all ROUGE metrics (Lin, 2004) as well  
 119           as on METEOR (Denkowski and Lavie, 2014)  
 120           of the CNN/Daily Mail dataset, achieving statistically significant improvements over previous  
 121           models that use complex long-encoder, copy, and coverage mechanisms (See et al., 2017). The  
 122           test-only DUC-2002 improvement also shows our  
 123           model’s better generalization than this strong ab-  
 124           stractive system. Moreover, our model’s training  
 125           is 4x and inference is more than 20x faster than  
 126           the previous state-of-the-art. The optional final  
 127           reranker gives further improvements while main-  
 128           taining a 7x speedup. In addition, we surpass the  
 129           popular lead-3 baseline on all ROUGE scores with  
 130           an abstractive model. This empirically justifies  
 131           that our RL-guided extractor has learnt sentence  
 132           saliency, rather than benefiting from simply copy-  
 133           ing longer sentences. We also show that our model  
 134           maintains the same level of fluency as a conven-  
 135           tional RNN-based model because the reward does  
 136           not leak to our abstractor’s word-level training.<sup>1</sup>

137           Overall, our contribution is three fold: First  
 138           we propose a novel sentence-level RL technique  
 139           for the well-known task of abstractive summariza-  
 140           tion, effectively utilizing the word-then-sentence  
 141           hierarchical structure without annotated matching  
 142           sentence-pairs between the document and ground  
 143           truth summary. Next, our model achieves the new  
 144           state-of-the-art on all metrics of multiple versions  
 145           of a popular summarization dataset (as well as a  
 146           test-only dataset) both extractively and abstrac-

147  
 148           <sup>1</sup>We plan to release our code, best trained models, as well  
 149           the generated summaries, to foster future research.

150           tively, without loss in language fluency. Finally,  
 151           our parallel decoding results in a significant 10-  
 152           20x speed-up over previous best neural abstractive  
 153           summarization system with even better accuracy.

## 2 Model

154           In this work, we consider the task of summariz-  
 155           ing a given long text document into several  
 156           (ordered) highlights, which are then combined  
 157           to form a multi-sentence summary. Formally,  
 158           given a training set of document-summary pairs  
 159            $\{x_i, y_i\}_{i=1}^N$ , our goal is to approximate the func-  
 160           tion  $h : X \rightarrow Y, X = \{x_i\}_{i=1}^N, Y = \{y_i\}_{i=1}^N$   
 161           such that  $h(x_i) = y_i, 1 \leq i \leq N$ . Furthermore,  
 162           we assume there exists an abstracting func-  
 163           tion  $g$  defined as:  $\forall s \in S_i, \exists d \in D_i$  such that  
 164            $g(d) = s, 1 \leq i \leq N$ , where  $S_i$  is the set of sum-  
 165           mary sentences in  $x_i$  and  $D_i$  the set of document  
 166           sentences in  $y_i$ . i.e., in any given pair of docu-  
 167           ment and summary, every summary sentence can  
 168           be produced from some document sentence. For  
 169           simplicity, we omit subscript  $i$  in the remainder  
 170           of the paper. Under this assumption, we can fur-  
 171           ther define another latent function  $f : X \rightarrow D^n$   
 172           that satisfies  $f(x) = \{d_t\}_{t=1}^n$  and  $y = h(x) = [g(d_1), g(d_2), \dots, g(d_n)]$ , where  $[,]$  denotes sen-  
 173           tence concatenation. This latent function  $f$  can be  
 174           seen as an extractor that chooses the salient (or-  
 175           dered) sentences in a given document for the ab-  
 176           stracting function  $g$  to rewrite. Our overall model  
 177           consists of these two submodules, the extractor  
 178           agent and the abstractor network, to approximate  
 179           the above-mentioned  $f$  and  $g$ , respectively.

### 2.1 Extractor Agent

180           The extractor agent is designed to model  $f$ , which  
 181           can be thought of as extracting salient sentences  
 182           from the document. We exploit a hierarchical neu-  
 183           ral model to learn the sentence representations of  
 184           the document and a ‘selection network’ to extract  
 185           sentences based on their representations.

#### 2.1.1 Hierarchical Sentence Representation

186           We use a temporal convolutional model (Kim,  
 187           2014) to compute  $r_j$ , the representation of each  
 188           individual sentence in the documents. The model  
 189           details are in the supplementary. To further incor-  
 190           porate global context of the document and capture  
 191           the long-range semantic dependency between sen-  
 192           tences, a bidirectional LSTM-RNN (Hochreiter  
 193           and Schmidhuber, 1997; Schuster et al., 1997) is  
 194           applied on the convolutional representation. This

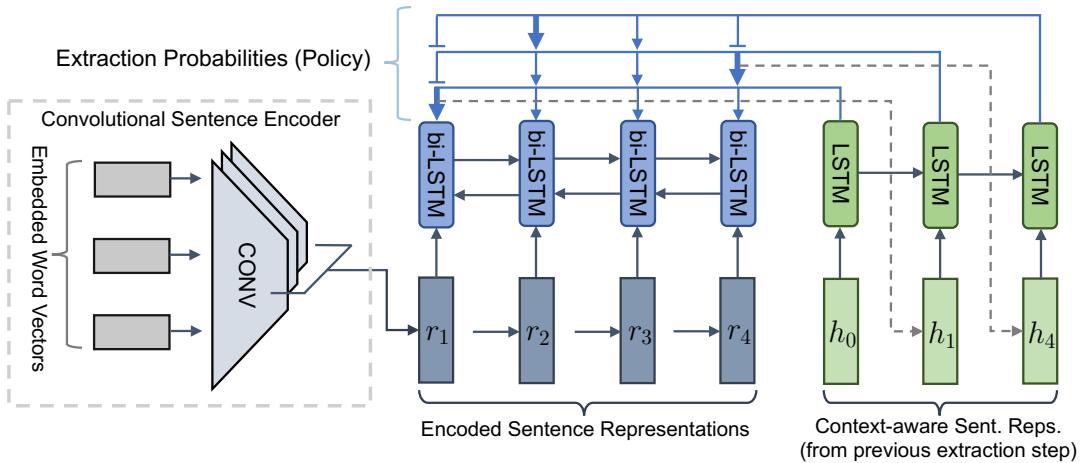


Figure 1: The sentence selection model. The convolutional encoder computes representation  $r_j$  for each sentence. The RNN encoder (blue) computes context-aware representation  $h_j$  and then the RNN decoder (green) selects sentence  $j_t$  at time step  $t$ . With  $j_t$  selected,  $h_{j_t}$  will be fed into the decoder at time  $t + 1$ .

enables learning a strong representation, denoted as  $h_j$  for the  $j$ -th sentence in the document, that takes into account the context of all previous and future sentences in the same document.

### 2.1.2 Sentence Selection

Next, to select the extracted sentences based on the above sentence representations, we add another **LSTM-RNN** to train a *Pointer Network* (Vinyals et al., 2015), to extract sentences recurrently. We calculate the extraction probability by:

$$u_j^t = \begin{cases} v_p^T \tanh(W_{p1}h_j + W_{p2}e_t) & \text{if } j_t \neq j_k \\ -\infty & \text{otherwise} \end{cases} \quad \forall k < t \quad (1)$$

$$P(j_t | j_1, \dots, j_{t-1}) = \text{softmax}(u^t) \quad (2)$$

where  $e_t$ 's are the output of the *glimpse* operation (Vinyals et al., 2016):

$$a_j^t = v_g^T \tanh(W_{g1}h_j + W_{g2}z_t) \quad (3)$$

$$\alpha^t = \text{softmax}(a^t) \quad (4)$$

$$e_t = \sum_j \alpha_j^t W_{g1}h_j \quad (5)$$

In Eqn. 3,  $z_t$  is the output of the added LSTM-RNN (shown in green in Fig. 1) which is referred as the *decoder*. All the  $W$ 's and  $v$ 's are trainable parameters. At each time step  $t$ , the decoder performs a 2-hop attention mechanism: It first attends to  $h_j$ 's to get a context vector  $e_t$  and then attends to  $h_j$ 's again for the extraction probabilities.<sup>2</sup> This

<sup>2</sup>In Eqn. 1, we zero the extraction prob. of already extracted sentences, a non-diff. op. only done in RL training.

model is essentially classifying all sentences of the document at each extraction step. An illustration of the whole extractor is shown in Fig. 1.

### 2.2 Abstractor Network

The abstractor network approximates  $g$ , which compresses and paraphrases an extracted document sentence to a concise summary sentence. We briefly discuss its architecture here and leave the details in the supplementary. We use the standard encoder-aligner-decoder model (Bahdanau et al., 2014; Luong et al., 2015), popular for several conditioned NL generation tasks. We add a simple copying mechanism<sup>3</sup> to our abstractor to help avoid generation of out-of-vocabulary(OOV) words (See et al., 2017), given the high overlap between the input and output of text summarization.

## 3 Learning

Given that our extractor performs a non-differentiable *hard extraction*, we apply standard policy gradient methods to bridge the back-propagation and form an end-to-end trainable (stochastic) computation graph. However, simply starting from randomly initialized network to train the whole model in an end-to-end fashion is infeasible. With randomly initialized networks, the extractor would often select sentences that are not relevant, so it would be difficult for the abstractor to learn to abstractively summarize. On the other hand, without a well-trained abstractor the extractor would get noisy reward, which

<sup>3</sup>We use *copy mechanism* (originally named *pointer-generator*) to avoid confusion with the *pointer network*.

leads to a bad estimate of policy gradient and a sub-optimal policy. We hence propose optimizing each sub-module separately using maximum-likelihood (ML) objectives: train the extractor to select salient sentences (fit  $f$ ) and the abstractor to generate shortened summary (fit  $g$ ). Finally, RL is applied to train the full model end-to-end (fit  $h$ ).

### 3.1 Maximum-Likelihood Training for Submodules

**Extractor Training:** In Sec. 2.1.2, we have formulated our sentence selection as classification. However, most of the summarization datasets are end-to-end document-summary pairs without extraction (saliency) labels of each sentence. Hence, we propose a simple similarity method to provide a ‘proxy’ target label for the extractor. Similar to the extractive model of Nallapati et al. (2017), for each ground-truth summary sentence, we find the most similar document sentence  $d_{j_t}$  by<sup>4</sup>:

$$j_t = \operatorname{argmax}_i (\text{ROUGE-L}_{\text{recall}}(d_i, s_t)) \quad (6)$$

Given these proxy training labels, the extractor is then trained to minimize the cross-entropy loss.

**Abstractor Training:** For the abstractor training, we create training pairs using each of document sentence extracted by Eqn. 6 above for each corresponding summary sentence. The network is trained as an usual sequence-to-sequence model to minimize the cross-entropy loss  $L(\theta_{abs}) = -\frac{1}{M} \sum_{m=1}^M \log P_{\theta_{abs}}(w_m | w_{1:m-1})$  of the decoder language model at each generation step, where  $\theta_{abs}$  is the set of trainable parameters of the abstractor and  $w_m$  the  $m^{\text{th}}$  generated word.

### 3.2 Reinforce-Guided Extraction

Here we explain how policy gradient techniques are applied to optimize the whole model. To make the extractor an RL agent, we can formulate a Markov Decision Process (MDP)<sup>5</sup>: at each extraction step  $t$ , the agent observes the current state  $c_t = (D, d_{j_{t-1}})$ , samples an action  $j_t \sim \pi_{\theta_a, \omega}(c_t, j) = P(j)$  from Eqn. 2 to extract a doc-

<sup>4</sup>Nallapati et al. (2017) also used a similar approach to obtain a proxy label for their extractive model; the main difference is that they chose sentences greedily to get the maximum summary-level ROUGE, whereas we choose sentences based on best individual matches of the sentence-level score to each ground-truth summary sentence (since we need to match exactly 1 document sentence for each GT summary sentence).

<sup>5</sup>Strictly speaking, this is a *Partially Observable Markov Decision Process* (POMDP). We approximate it by MDP by assuming that the RNN hidden state contains all past info.

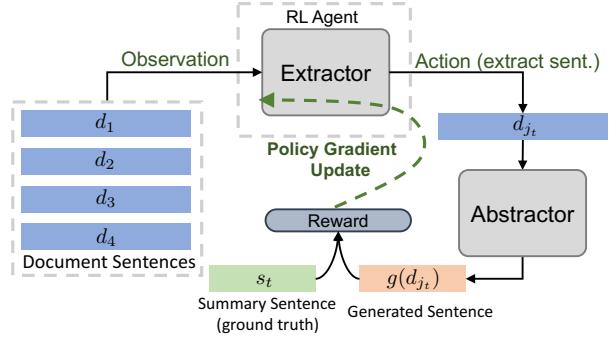


Figure 2: Reinforced training of the extractor (for one extraction step). For simplicity, the critic network is not shown. Note that all  $d$ 's and  $s_t$  are raw sentences, not some vector representations.

ument sentence and receive a reward<sup>6</sup>

$$r(t+1) = \text{ROUGE-L}_{F_1}(g(d_{j_t}), s_t) \quad (7)$$

after the abstractor summarizes the extracted sentence  $d_{j_t}$ . We denote the trainable parameters of the extractor agent by  $\theta = \{\theta_a, \omega\}$  for the decoder and hierarchical encoder respectively. We can then train the extractor with policy-based RL.

The vanilla policy gradient algorithm, REINFORCE (Williams, 1992), is known for high variance. To mitigate this problem we add a critic network with trainable parameters  $\theta_c$  to predict the state-value function  $V^{\pi_{\theta_a, \omega}}(c)$ . The predicted value of critic  $b_{\theta_c, \omega}(c)$  is called the ‘baseline’, which is used to estimate the advantage function:  $A^{\pi_{\theta}}(c, j) = Q^{\pi_{\theta_a, \omega}}(c, j) - V^{\pi_{\theta_a, \omega}}(c)$  because the total return  $R_t$  is an estimate of action-value function  $Q(c_t, j_t)$ . Instead of maximizing  $Q(c_t, j_t)$  as done in REINFORCE, we maximize  $A^{\pi_{\theta}}(c, j)$  with the following policy gradient:

$$\nabla_{\theta_a, \omega} J(\theta_a, \omega) = \mathbb{E}[\nabla_{\theta_a, \omega} \log \pi_{\theta}(c, j) A^{\pi_{\theta}}(c, j)] \quad (8)$$

And the critic is trained to minimize the square loss:  $L_c(\theta_c, \omega) = (b_{\theta_c, \omega}(c_t) - R_t)^2$ . This is known as the *Advantage Actor-Critic* (A2C), a synchronous variant of A3C (Mnih et al., 2016). For more A2C details, please refer to the supp.

Intuitively, our RL training works as follow: If the extractor chooses a good sentence, after the abstractor rewrites it the ROUGE match would be high and thus the action is encouraged. If a bad sentence is chosen, though the abstractor

<sup>6</sup>In Eqn. 6 we use ROUGE-recall because we want the extracted sentence to contain as much information as possible for rewriting. Nevertheless, for Eqn. 7, ROUGE- $F_1$  is more suitable because the abstractor  $g$  is supposed to rewrite the extracted sentence  $d$  to be as concise as the ground truth  $s$ .

still produce a compressed version of it, the summary would not match the ground truth and the low ROUGE score discourages this action. Our RL with a sentence-level agent is a novel attempt in neural summarization. We use RL as a saliency guide without altering the language model, while previous work applied RL on the word-level, which could be prone to gaming the metric at the cost of language fluency.<sup>7</sup>

### 3.3 Repetition-Avoiding Reranking

Existing abstractive summarization systems on long documents suffer from generating repeating and redundant words and phrases. To mitigate this issue, See et al. (2017) proposes the coverage mechanism and Paulus et al. (2017) incorporates tri-gram avoidance during beam-search at test-time. Our model without these already performs well because the summary sentences are generated from different document sentences, which naturally avoid redundancy. However, we do get a small further boost to the summary quality by removing a few ‘across-sentence’ repetitions, via a simple reranking strategy: At sentence-level, we apply the same beam-search tri-gram avoidance (Paulus et al., 2017). We keep all  $k$  sentence candidates generated by beam search, where  $k$  is the size of the beam. Next, we then rerank all  $k^n$  combinations of the  $n$  generated summary sentence beams. The summaries are reranked by the number of repeated N-grams, the smaller the better. We also apply the diverse decoding algorithm described in Li et al. (2016) (which has almost no computation overhead) so as to get the above approach to produce useful diverse reranking lists. We show how much the redundancy affects the summarization task in Sec. 6.2.1.

## 4 Related Work

Early summarization works mostly focused on extractive and compression based methods (Jing and McKeown, 2000; Knight and Marcu, 2000; Clarke and Lapata, 2010; Berg-Kirkpatrick et al., 2011;

<sup>7</sup>During this RL training of the extractor, we keep the abstractor parameters fixed. Because the input sentences for the abstractor are extracted by an intermediate stochastic policy of the extractor, it is impossible to find the correct target summary for the abstractor to fit  $g$  with ML objective. Though it is possible to optimize the abstractor with RL (Paulus et al., 2017), in our experiments we found that this does not improve the overall ROUGE, most likely because this RL optimizes at a sentence-level and can add across-sentence redundancy. We achieve SotA results without this abstractor-level RL.

Filippova et al., 2015). Recent large-sized corpora attracted neural methods for abstractive summarization (Rush et al., 2015; Chopra et al., 2016). Some of the recent success in neural abstractive models include hierarchical attention (Nallapati et al., 2016), coverage (Suzuki and Nagata, 2016; Chen et al., 2016; See et al., 2017), RL based metric optimization (Paulus et al., 2017), graph-based attention (Tan et al., 2017b), and the copy mechanism (Miao and Blunsom, 2016; Gu et al., 2016a; See et al., 2017), which combines extractive and abstractive methods in a soft way using attention.

Our model shares some high-level intuition with extract-then-compress methods. Earlier attempts in this paradigm used Hidden Markov Models and rule-based systems (Jing and McKeown, 2000), statistical models based on parse trees (Knight and Marcu, 2000), and integer linear programming based methods (Martins and Smith, 2009; Gillick and Favre, 2009; Clarke and Lapata, 2010; Berg-Kirkpatrick et al., 2011). Recent approaches investigated discourse structure (Louis et al., 2010; Hirao et al., 2013; Kikuchi et al., 2014; Wang et al., 2015), graph cuts (Qian and Liu, 2013), and parse trees (Li et al., 2014; Bing et al., 2015). For neural models, Cheng and Lapata (2016) used a second neural net to select words from an extractor’s output. Our abstractor does not merely ‘compress’ the sentences but generatively produce novel words. Moreover, our RL bridges the extractor and the abstractor for end-to-end training.

Reinforcement learning has been used to optimize the non-differential metrics of language generation and to mitigate exposure bias (Ranzato et al., 2015; Bahdanau et al., 2016). Henß et al. (2015) use Q-learning based RL for extractive summarization. Paulus et al. (2017) use RL policy gradient methods for abstractive summarization, utilizing sequence-level metric rewards with curriculum learning (Ranzato et al., 2015) or weighted ML+RL mixed loss (Paulus et al., 2017) for stability and language fluency. We use sentence-level rewards to optimize the extractor while keeping our ML trained abstractor decoder fixed, so as to achieve the best of both worlds.

Training a neural network to use another fixed network has been investigated in machine translation for better decoding (Gu et al., 2017) and real-time translation (Gu et al., 2016b). They used a fixed pretrained translator (similar to our abstractor) and applied policy gradient techniques

to train another task-specific network. In question answering (QA), Choi et al. (2017) extract sentences and then generate the answer from their vector representations. We instead use the raw selected sentences to employ the copy mechanism; also, we produce multi-sentence summaries instead of just one answer. Another recent work attempted a new coarse-to-fine attention approach on summarization (Ling and Rush, 2017) and found desired sharp focus properties for scaling to larger inputs (though without metric improvements). Lastly, our extractor is inspired by a recent attempt on solving the traveling salesman problem using pointer networks and RL (Bello et al., 2017).

Finally, there are some loosely-related recent works: Zhou et al. (2017) proposed *selective gate* to improve the attention in abstractive summarization. Tan et al. (2017a) uses an *extract-then-synthesis* approach on QA, where an extraction model predicts the important spans in the passage and then another synthesis model generates the final answer. Swayamdipta et al. (2017) attempted cascaded non-recurrent small networks on extractive QA, resulting a scalable, parallelizable model. Fan et al. (2017) added controlling parameters to allow summarization to user preferences of length, style, etc. However, none of these used RL to bridge non-differentiability of neural models.

## 5 Experimental Setup

We use the CNN/Daily Mail dataset first proposed by Hermann et al. (2015) and then later modified for summarization by Nallapati et al. (2017). Because there are two versions of the dataset, original text and entity anonymized, we show results on both versions of the dataset for a fair comparison to prior works. The experiment runs training and evaluation for each version separately. Despite the fact that the 2 versions have been considered separately by the summarization community as 2 different datasets, we use same hyper-parameter values for both dataset versions to show the generalization of our model. We also show improvements on the DUC-2002 dataset in a *test-only* setup. For more details, please refer to the supplementary.

### 5.1 Evaluation Metrics

For all the datasets, we evaluate standard ROUGE-1, ROUGE-2, and ROUGE-L (Lin, 2004) on full-length  $F_1$  (with stemming) following previous works (Nallapati et al., 2017; See et al., 2017;

Paulus et al., 2017). Following See et al. (2017), we also evaluate on METEOR (Denkowski and Lavie, 2014) for a more thorough analysis.

## 5.2 Modular Extractive vs. Abstractive

Our hybrid extractive-abstractive approach is capable of both extractive and abstractive (i.e., rewriting every sentence) summarization. The extractor alone performs extractive summarization. To investigate the effect of the recurrent extractor (rnn-ext), we implement a feed-forward extractive baseline ff-ext(details in supp.). It is also possible to apply RL to extractor without using the abstractor (rnn-ext + RL).<sup>8</sup> Benefiting from the high modularity of our model, we can make our summarization system abstractive by simply applying the abstractor on the extracted sentences. Our abstractor rewrites each sentence and generates novel words from a large vocabulary, and hence every word in our overall summary is generated from scratch; making our full model categorized into the abstractive paradigm.<sup>9</sup> We run experiments on separately trained extractor/abstractor (ff-ext + abs, rnn-ext + abs) and the reinforced full model (rnn-ext + abs + RL) as well as the final reranking version (rnn-ext + abs + RL + rerank).

## 6 Results

For easier comparison, we show separate tables for the original-text vs. anonymized versions (with corresponding previous works in each table) – Table 1 and Table 2, respectively. Overall, our model achieves strong improvements and the new state-of-the-art on both extractive and abstractive settings for both versions of the CNN/DM dataset (with some comparable results on the anonymized version). Moreover, Table 3 shows the generalization of our abstractive system to an out-of-domain test-only setup (DUC-2002), where our model achieves better scores than See et al. (2017).

### 6.1 Extractive Summarization

In the extractive paradigm, we compare our model with the extractive model from Nallapati et al. (2017) and a strong lead-3 baseline. For producing our summary, we simply concatenate the extracted sentences from the extractors. From Table 2 and

<sup>8</sup>In this case the abstractor function  $g(d) = d$ .

<sup>9</sup>Note that the CNN/DM dataset does *not* include any extraction label which could have benefited our abstractive full model over previous abstractive models.

Models	ROUGE-1	ROUGE-2	ROUGE-L	METEOR
Extractive Results				
lead-3 (See et al., 2017)	40.34	17.70	36.57	22.21
ff-ext	<b>40.63</b>	<b>18.35</b>	<b>36.82</b>	<b>22.91</b>
rnn-ext	40.17	18.11	36.41	22.81
rnn-ext + RL	<b>41.47</b>	<b>18.72</b>	<b>37.76</b>	22.35
Abstractive Results				
See et al. (2017) (w/o coverage)	36.44	15.66	33.42	16.65
See et al. (2017)	39.53	<b>17.28</b>	36.38	18.72
Fan et al. (2017) (controlled)	39.75	<b>17.29</b>	36.54	-
ff-ext + abs	39.30	17.02	36.93	20.05
rnn-ext + abs	38.38	16.12	36.04	19.39
rnn-ext + abs + RL	<b>40.04</b>	<b>17.61</b>	<b>37.59</b>	<b>21.00</b>
rnn-ext + abs + RL + rerank	<b>40.88</b>	<b>17.80</b>	<b>38.54</b>	20.38

Table 1: Results on the original, non-anonymized CNN/Daily Mail dataset. Adding RL gives statistically significant improvements for all metrics over non-RL rnn-ext models (and over the state-of-the-art See et al. (2017)) in both extractive and abstractive settings with  $p < 0.01$ . Adding the extra reranking stage yields statistically significant better results in terms of all ROUGE metrics with  $p < 0.01$ .

Models	R-1	R-2	R-L
Extractive Results			
lead-3 (Nallapati et al., 2017)	39.20	15.70	35.5
Nallapati et al. (2017)	39.60	16.20	35.30
ff-ext	39.51	<b>16.85</b>	35.80
rnn-ext	38.97	16.65	35.32
rnn-ext + RL	<b>40.13</b>	16.58	<b>36.47</b>
Abstractive Results			
Nallapati et al. (2016)	35.46	13.30	32.65
Fan et al. (2017) (controlled)	38.68	15.40	35.47
Paulus et al. (2017) (ML)	38.30	14.81	35.49
Paulus et al. (2017) (RL+ML)	<b>39.87</b>	15.82	36.90
ff-ext + abs	38.73	15.70	36.33
rnn-ext + abs	37.58	14.68	35.24
rnn-ext + abs + RL	38.80	15.66	36.37
rnn-ext + abs + RL + rerank	39.66	<b>15.85</b>	<b>37.34</b>

Table 2: ROUGE for anonymized CNN/DM.

Table 1, we can see that our feed-forward extractor out-performs the lead-3 baseline, empirically showing that our hierarchical sentence encoding model is capable of extracting salient sentences.<sup>10</sup> The reinforced extractor performs the best, because of the ability to get the summary-level reward and the reduced train-test mismatch of feeding the previous extraction decision. The improvement over lead-3 is consistent across both tables. It also outperforms the previous best neural extractive model (Nallapati et al., 2017).

## 6.2 Abstractive Summarization

After applying the abstractor, the ff-ext based model still out-performs the rnn-ext model. Both combined models exceeds the pointer-generator model (See et al., 2017) without coverage by a large margin for all metrics, showing the effec-

<sup>10</sup>We suspect ff-ext outperforms rnn-ext because feed-forward model is easier to optimize in nature and the evaluations (n-gram matching) do not consider sentence ordering.

Models	R-1	R-2	R-L
See et al. (2017)	37.22	<b>15.78</b>	33.90
rnn-ext + abs + RL	39.46	<b>17.34</b>	36.72

Table 3: Generalization to DUC-2002.

tiveness of our 2-step hierarchical approach: our method naturally avoids repetition by extracting multiple sentences with different keypoints.

Moreover, after applying reinforcement learning, our model performs better than the best model of See et al. (2017) and the best ML trained model of Paulus et al. (2017). Our reinforced model outperforms the ML trained rnn-ext + abs baseline with statistical significance of  $p < 0.01$  on all metrics for both version of the dataset, indicating the effectiveness of the RL training. Also, rnn-ext + abs + RL is statistically significant better than See et al. (2017) for all metrics with  $p < 0.01$ .<sup>11</sup> In the supplementary, we show the learning curve of our RL training, where the average reward goes up quickly after the extractor learned the End-of-Extract action and then stabilizes. For all the above models, we use standard greedy decoding and find that it performs well.<sup>12</sup>

### 6.2.1 Reranking and Redundancy

Although the extract-then-abstract approach inherently will not generate repeating sentences like

<sup>11</sup>We calculate statistical significance based on the bootstrap test (Noreen, 1989; Efron and Tibshirani, 1994) with 100K samples. Output of Paulus et al. (2017) is not available so we couldn't test for statistical significance there.

<sup>12</sup>In our MDP formulation, we cannot apply RL on ff-ext due to its historyless nature. It is, however, possible to apply RL naively. We did not do so because we do not want to rely on the ordering assumption (Eqn. 5 in the supp.), though the current metrics are agnostic to sentence ordering. Our rnn-ext therefore is more generalizable in the sense of less assumption and will be easier to adapt to arbitrary metrics.

Models	Speed	
	total time (hr)	words / sec
(See et al., 2017)	12.9	14.8
rnn-ext + abs + RL	0.68	361.3
rnn-ext + abs + RL + rerank	2.00 (1.46 +0.54)	109.8

Table 4: Speed comparison with See et al. (2017).

other neural-decoders do, there might still be across-sentence redundancy because the abstractor is not aware of other extracted sentences when decoding one. Hence, we incorporate an optional reranking strategy described in Sec. 3.3. The improved ROUGE scores indicate that this successfully removes some remaining redundancies and hence produces more concise summaries. Our best abstractive model (rnn-ext + abs + RL + rerank) is clearly superior than the one of See et al. (2017). We are comparable on R-1 and R-2 but a 0.4 point improvement on R-L w.r.t. Paulus et al. (2017).<sup>13</sup>

Several previous works have pointed out that extractive baselines are very difficult to beat (in terms of ROUGE) by an abstractive system (See et al., 2017; Nallapati et al., 2017). Note that our best model is the first abstractive model to report better ROUGE scores than the lead-3 baseline on the original-text CNN/DM dataset. Our extractive experiment serves as a complementary analysis on the effect of RL with extractive systems.

### 6.3 Speed Comparison

Our two-stage extractive-abstractive hybrid model is not only the SotA on summary quality metrics, but more importantly also gives a significant speed-up in both train and test time over a strong neural abstractive system (See et al., 2017).<sup>14</sup>

Our full model is composed of a extremely fast extractor and a parallelizable abstractor, where the computation bottleneck is on the abstractor, which has to generate summaries with a large vocabulary from scratch.<sup>15</sup> The main advantage of our abstractor at decoding time is that we can first compute all the extracted sentences for the document, and then abstract every sentence concurrently (in

<sup>13</sup>We do not list the scores of their pure RL model because they discussed its bad readability. Moreover, we also report new results from the recent preprint of Fan et al. (2017)—our model achieves better scores over this.

<sup>14</sup>The only publicly available code with a pretrained model for neural summarization which we can test the speed.

<sup>15</sup>The time needed for extractor is negligible w.r.t. the abstractor because it does not require large matrix multiplication like generative model does for every word. Moreover, with convolutional encoder (very efficient with GPU) at word-level made parallelizable by the hierarchical rnn-ext, our model is scalable for very long documents.

parallel) to generate the overall summary. In Table 4, we show the substantial test-time speed-up of our model compared to See et al. (2017).<sup>16</sup> We calculate the total decoding time for producing all summaries for the test set.<sup>17</sup> Due to the fact that the main test-time speed bottleneck of RNN language generation model is that the model is constrained to generate one word at a time, the total decoding time is dependent on the number of total words generated; we hence also report the decoded words per second for a fair comparison. Our model without reranking is extremely fast. From Table 4 we can see that we achieve a speed up of 18x in time and 24x in word generation rate. Even after adding the (optional) reranker, we still maintain a 6-7x speed-up (and hence a user can choose to use the reranking component depending on their downstream application’s speed requirements).

### 6.4 Qualitative Analysis

Overall we observe that extractive methods consistently get higher ROUGE  $F_1$  scores over abstractive models. However, abstractive models generate shorter sentences that are more concise. We show examples on how our best model selects sentences and then rewrites them. In the supplementary Table 2, we can see that the abstractor does compress the sentence to be more concise while keeping correct facts. Also, the reranker helps removing redundant information in the second and third summary sentence. We observe that when rewriting longer text, the abstractor would have many facts to choose and this is when the reranker helped avoid redundancy across sentences. We do not observe any fluency issue (Paulus et al., 2017). In the supplementary we empirically show that our models’ summaries achieve better/comparable perplexity w.r.t. the ground truth.

## 7 Conclusion

We proposed a novel sentence-level RL model for abstractive summarization, which makes the model aware of the word-sentence hierarchy. Our model achieved the new state-of-the-art on both CNN/DM versions as well a better generalization on test-only DUC-2002, along with a significant speed-up in training and decoding.

<sup>16</sup>For details of training speed-up, please see the supp.

<sup>17</sup>We time the model of See et al. (2017) using beam size of 4 (used for their best-reported scores). Without beam-search, it gets significantly worse ROUGE of 36.62, 15.12, 34.08, so we don’t compare speed-ups w.r.t. that version.

## 800 References

- 801 Dzmitry Bahdanau, Philemon Brakel, Kelvin Xu,  
802 Anirudh Goyal, Ryan Lowe, Joelle Pineau, Aaron C.  
803 Courville, and Yoshua Bengio. 2016. **An actor-**  
804 **critic algorithm for sequence prediction.** *CoRR*,  
805 abs/1607.07086.
- 806 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Ben-  
807 gio. 2014. Neural machine translation by jointly  
808 learning to align and translate. *arXiv e-prints*,  
809 abs/1409.0473.
- 810 Michele Banko, Vibhu O. Mittal, and Michael J. Wit-  
811 brock. 2000. **Headline generation based on statis-**  
812 **tical translation.** In *Proceedings of the 38th An-*  
813 *nual Meeting on Association for Computational Lin-*  
814 *guistics*, ACL '00, pages 318–325, Stroudsburg, PA,  
USA. Association for Computational Linguistics.
- 815 Irwan Bello, Hieu Pham, Quoc V. Le, Mohammad  
816 Norouzi, and Samy Bengio. 2017. **Neural combi-**  
817 **natorial optimization with reinforcement learning.**
- 818 Taylor Berg-Kirkpatrick, Dan Gillick, and Dan Klein.  
819 2011. **Jointly learning to extract and compress.** In  
820 *Proceedings of the 49th Annual Meeting of the Asso-*  
821 *ciation for Computational Linguistics: Human Lan-*  
822 *guage Technologies - Volume 1*, HLT '11, pages  
481–490, Stroudsburg, PA, USA. Association for  
823 Computational Linguistics.
- 824 Lidong Bing, Piji Li, Yi Liao, Wai Lam, Weiwei  
825 Guo, and Rebecca J. Passonneau. 2015. Abstractive  
826 multi-document summarization via phrase selection  
827 and merging. In *ACL (1)*, pages 1587–1597. The  
828 Association for Computer Linguistics.
- 829 Qian Chen, Xiaodan Zhu, Zhenhua Ling, Si Wei, and  
830 Hui Jiang. 2016. Distraction-based neural networks  
831 for modeling documents. In *IJCAI*.
- 832 Jianpeng Cheng and Mirella Lapata. 2016. **Neural**  
833 **summarization by extracting sentences and words.**  
834 In *Proceedings of the 54th Annual Meeting of the*  
835 *Association for Computational Linguistics (Volume*  
836 *1: Long Papers)*, pages 484–494, Berlin, Germany.  
Association for Computational Linguistics.
- 837 Eunsol Choi, Daniel Hewlett, Jakob Uszkoreit, Illia  
838 Polosukhin, Alexandre Lacoste, and Jonathan Be-  
839 rant. 2017. **Coarse-to-fine question answering for**  
840 **long documents.** In *Proceedings of the 55th Annual*  
841 *Meeting of the Association for Computational Lin-*  
842 *guistics (Volume 1: Long Papers)*, pages 209–220.  
Association for Computational Linguistics.
- 843 Sumit Chopra, Michael Auli, and Alexander M. Rush.  
2016. **Abstractive sentence summarization with at-**  
844 **tentive recurrent neural networks.** In *Proceedings of*  
845 *the 2016 Conference of the North American Chapter*  
846 *of the Association for Computational Linguistics:*  
847 *Human Language Technologies*, pages 93–98, San  
Diego, California. Association for Computational  
848 Linguistics.
- James Clarke and Mirella Lapata. 2010. Discourse  
constraints for document compression. *Computational  
Linguistics*, 36(3):411–441.
- Michael Denkowski and Alon Lavie. 2014. Meteor  
universal: Language specific translation evaluation  
for any target language. In *Proceedings of the EACL  
2014 Workshop on Statistical Machine Translation*.
- Bradley Efron and Robert J Tibshirani. 1994. *An intro-  
duction to the bootstrap.* CRC press.
- Angela Fan, David Grangier, and Michael Auli. 2017.  
**Controllable abstractive summarization.** *CoRR*,  
abs/1711.05217.
- Katja Filippova, Enrique Alfonseca, Carlos Col-  
menares, Lukasz Kaiser, and Oriol Vinyals. 2015.  
Sentence compression by deletion with lstms. In  
*Proceedings of the 2015 Conference on Empirical  
Methods in Natural Language Processing (EMNLP'15)*.
- Dan Gillick and Benoit Favre. 2009. **A scalable global  
model for summarization.** In *Proceedings of the  
Workshop on Integer Linear Programming for Nat-  
ural Langauge Processing*, ILP '09, pages 10–18,  
Stroudsburg, PA, USA. Association for Compu-  
tational Linguistics.
- Jiatao Gu, Kyunghyun Cho, and Victor O. K. Li. 2017.  
**Trainable greedy decoding for neural machine trans-  
lation.** *CoRR*, abs/1702.02429.
- Jiatao Gu, Zhengdong Lu, Hang Li, and Victor O.K.  
Li. 2016a. Incorporating copying mechanism in  
sequence-to-sequence learning. In *Proceedings of  
the 54th Annual Meeting of the Association for Com-  
putational Linguistics (Volume 1: Long Papers)*,  
pages 1631–1640, Berlin, Germany. Association for  
Computational Linguistics.
- Jiatao Gu, Graham Neubig, Kyunghyun Cho, and Vic-  
tor O. K. Li. 2016b. **Learning to translate in  
real-time with neural machine translation.** *CoRR*,  
abs/1610.00388.
- Sebastian Henß, Margot Mieskes, and Iryna Gurevych.  
2015. A reinforcement learning approach for adap-  
tive single- and multi-document summarization. In  
*International Conference of the German Society for  
Computational Linguistics and Language Technol-  
ogy (GSCL-2015)*, pages 3–12.
- Karl Moritz Hermann, Tomáš Kočiský, Edward  
Grefenstette, Lasse Espeholt, Will Kay, Mustafa Su-  
leyman, and Phil Blunsom. 2015. **Teaching ma-  
chines to read and comprehend.** In *Advances in Neu-  
ral Information Processing Systems (NIPS)*.
- Tsutomu Hirao, Yasuhisa Yoshida, Masaaki Nishino,  
Norihito Yasuda, and Masaaki Nagata. 2013.  
**Single-document summarization as a tree knapsack  
problem.** In *Proceedings of the 2013 Conference on  
Empirical Methods in Natural Language Process-  
ing*, pages 1515–1520, Seattle, Washington, USA.  
Association for Computational Linguistics.

- 900 Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long  
901 short-term memory. *Neural Comput.*, 9(9):1735–  
902 1780.
- 903 Hongyan Jing and Kathleen R. McKeown. 2000. **Cut**  
904 **and paste based text summarization**. In *Proceed-  
905 ings of the 1st North American Chapter of the As-  
906 sociation for Computational Linguistics Conference,  
907 NAACL 2000*, pages 178–185, Stroudsburg, PA,  
908 USA. Association for Computational Linguistics.
- 909 Yuta Kikuchi, Tsutomu Hirao, Hiroya Takamura, Man-  
910 abu Okumura, and Masaaki Nagata. 2014. **Single**  
911 **document summarization based on nested tree struc-  
912 ture**. In *Proceedings of the 52nd Annual Meeting of  
913 the Association for Computational Linguistics (Vol-  
914 ume 2: Short Papers)*, pages 315–320, Baltimore,  
Maryland. Association for Computational Linguis-  
tics.
- 915 Yoon Kim. 2014. Convolutional neural networks  
916 for sentence classification. In *Proceedings of the  
917 2014 Conference on Empirical Methods in Natural  
918 Language Processing (EMNLP)*, pages 1746–1751,  
Doha, Qatar. Association for Computational Lin-  
guistics.
- 921 Kevin Knight and Daniel Marcu. 2000. **Statistics-  
922 based summarization - step one: Sentence compres-  
923 sion**. In *Proceedings of the Seventeenth National  
924 Conference on Artificial Intelligence and Twelfth  
925 Conference on Innovative Applications of Artificial  
Intelligence*, pages 703–710. AAAI Press.
- 926 Chen Li, Yang Liu, Fei Liu, Lin Zhao, and Fuliang  
927 Weng. 2014. **Improving multi-documents summa-  
928 rization by sentence compression based on expanded  
929 constituent parse trees**. In *Proceedings of the 2014  
930 Conference on Empirical Methods in Natural Lan-  
931 guage Processing (EMNLP)*, pages 691–701. Asso-  
ciation for Computational Linguistics.
- 932 Jiwei Li, Will Monroe, and Dan Jurafsky. 2016. **A sim-  
933 ple, fast diverse decoding algorithm for neural gen-  
934 eration**. *CoRR*, abs/1611.08562.
- 935 Chin-Yew Lin. 2004. **Rouge: A package for automatic  
936 evaluation of summaries**. In *Text Summarization  
937 Branches Out: Proceedings of the ACL-04 Work-  
938 shop*, pages 74–81, Barcelona, Spain. Association  
939 for Computational Linguistics.
- 940 Jeffrey Ling and Alexander Rush. 2017. **Coarse-to-fine  
941 attention models for document summarization**. In  
942 *Proceedings of the Workshop on New Frontiers in  
943 Summarization*, pages 33–42. Association for Com-  
putational Linguistics.
- 944 Annie Louis, Aravind Joshi, and Ani Nenkova. 2010.  
945 **Discourse indicators for content selection in summa-  
946 rization**. In *Proceedings of the 11th Annual Meeting  
947 of the Special Interest Group on Discourse and Dia-  
logue*, SIGDIAL ’10, pages 147–156, Stroudsburg,  
PA, USA. Association for Computational Linguis-  
tics.
- 948 Minh-Thang Luong, Hieu Pham, and Christopher D.  
949 Manning. 2015. Effective approaches to attention-  
based neural machine translation. In *Empirical  
Methods in Natural Language Processing (EMNLP)*,  
pages 1412–1421, Lisbon, Portugal. Association for  
Computational Linguistics.
- 950 André F. T. Martins and Noah A. Smith. 2009. **Sum-  
951 marization with a joint model for sentence extraction  
952 and compression**. In *Proceedings of the Workshop  
953 on Integer Linear Programming for Natural Lan-  
954 guage Processing*, ILP ’09, pages 1–9, Stroudsburg,  
PA, USA. Association for Computational Linguis-  
tics.
- 955 Yishu Miao and Phil Blunsom. 2016. Language as a  
956 latent variable: Discrete generative models for sen-  
957 tence compression. *CoRR*, abs/1609.07317.
- 958 Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi  
959 Mirza, Alex Graves, Timothy Lillicrap, Tim Harley,  
David Silver, and Koray Kavukcuoglu. 2016. **Asyn-  
960 chronous methods for deep reinforcement learning**.  
In *Proceedings of The 33rd International Confer-  
961 ence on Machine Learning*, volume 48 of *Proceed-  
962 ings of Machine Learning Research*, pages 1928–  
1937, New York, New York, USA. PMLR.
- 963 Ramesh Nallapati, Feifei Zhai, and Bowen Zhou. 2017.  
964 **Summarunner: A recurrent neural network based se-  
965 quence model for extractive summarization of doc-  
966 uments**. In *AAAI Conference on Artificial Intelli-  
967 gence*.
- 968 Ramesh Nallapati, Bowen Zhou, Cicero Nogueira dos  
969 santos, Caglar Gulcehre, and Bing Xiang. 2016.  
970 **Abstractive text summarization using sequence-to-  
971 sequence rnns and beyond**. In *CoNLL*.
- 972 Eric W Noreen. 1989. *Computer-intensive methods for  
973 testing hypotheses*. Wiley New York.
- 974 Romain Paulus, Caiming Xiong, and Richard Socher.  
2017. **A deep reinforced model for abstractive sum-  
975 marization**. *arXiv preprint*, abs/1705.04304.
- 976 Xian Qian and Yang Liu. 2013. **Fast joint compres-  
977 sion and summarization via graph cuts**. In *Proceed-  
978 ings of the 2013 Conference on Empirical Methods  
979 in Natural Language Processing*, pages 1492–1502,  
Seattle, Washington, USA. Association for Compu-  
tational Linguistics.
- 980 Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli,  
981 and Wojciech Zaremba. 2015. **Sequence level  
982 training with recurrent neural networks**. *CoRR*,  
abs/1511.06732.
- 983 Alexander M. Rush, Sumit Chopra, and Jason Weston.  
2015. **A neural attention model for abstractive sen-  
984 tence summarization**. In *Proceedings of the 2015  
985 Conference on Empirical Methods in Natural Lan-  
986 guage Processing*, pages 379–389, Lisbon, Portugal.  
987 Association for Computational Linguistics.
- 988 989 990 991 992 993 994 995 996 997 998 999

- 1000 Mike Schuster, Kuldip K. Paliwal, and A. General. 1050  
 1001 1997. Bidirectional recurrent neural networks. 1051  
 1002 *IEEE Transactions on Signal Processing*. 1052
- 1003 Abigail See, Peter J. Liu, and Christopher D. Manning. 1053  
 1004 2017. Get to the point: Summarization with pointer- 1054  
 1005 generator networks. In *Proceedings of the 55th An- 1055  
 1006 nual Meeting of the Association for Computational 1056  
 1007 Linguistics (Volume 1: Long Papers)*, pages 1073– 1057  
 1008 1083. Association for Computational Linguistics. 1058
- 1009 Jun Suzuki and Masaaki Nagata. 2016. Rnn-based 1059  
 1010 encoder-decoder approach with word frequency es- 1060  
 1011 timation. In *EACL*. 1060
- 1012 Swabha Swayamdipta, Ankur P. Parikh, and Tom 1061  
 1013 Kwiatkowski. 2017. Multi-mention learning 1062  
 1014 for reading comprehension with neural cascades. 1063  
*CoRR*, abs/1711.00894. 1064
- 1015 Chuanqi Tan, Furu Wei, Nan Yang, Weifeng Lv, and 1065  
 1016 Ming Zhou. 2017a. S-net: From answer extraction 1066  
 1017 to answer generation for machine reading compre- 1067  
 1018 hension. *CoRR*, abs/1706.04815. 1068
- 1019 Jiwei Tan, Xiaojun Wan, and Jianguo Xiao. 2017b. 1069  
 1020 Abstractive document summarization with a graph- 1070  
 1021 based attentional neural model. In *ACL*. 1071
- 1022 Oriol Vinyals, Samy Bengio, and Manjunath Kudlur. 1072  
 1023 2016. Order matters: Sequence to sequence for sets. 1073  
 1024 In *International Conference on Learning Represen- 1074  
 1025 tations (ICLR)*. 1075
- 1026 Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. 1076  
 1027 2015. Pointer networks. In C. Cortes, N. D. 1077  
 1028 Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, 1078  
 1029 editors, *Advances in Neural Information Processing 1079  
 1030 Systems 28*, pages 2692–2700. Curran Associates, 1080  
 Inc. 1081
- 1031 Xun Wang, Yasuhisa Yoshida, Tsutomu Hirao, Kat- 1082  
 1032 suhito Sudoh, and Masaaki Nagata. 2015. Sum- 1083  
 1033 marization based on task-oriented discourse parsing. 1084  
*IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, 1085  
 1034 23(8):1358–1367. 1086
- 1035 Ronald J. Williams. 1992. Simple statistical gradient- 1087  
 1036 following algorithms for connectionist reinforce- 1088  
 1037 ment learning. *Mach. Learn.*, 8(3-4):229–256. 1089
- 1038 David Zajic, Bonnie Dorr, and Richard Schwartz. 2004. 1090  
 1039 Bbn/umd at duc-2004: Topiary. In *HLT-NAACL 1091  
 1040 2004 Document Understanding Workshop*, pages 112–119, Boston, Massachusetts. 1092
- 1041 Qingyu Zhou, Nan Yang, Furu Wei, and Ming Zhou. 1093  
 1042 2017. Selective encoding for abstractive sentence 1094  
 1043 summarization. In *Proceedings of the 55th An- 1095  
 1044 nual Meeting of the Association for Computational 1096  
 1045 Linguistics (Volume 1: Long Papers)*, pages 1095– 1097  
 1046 1104. Association for Computational Linguistics. 1098
- 1047 1048 1049 1099