

# Improving Multi-Document Summarization through Referenced Flexible Extraction with Credit-Awareness

Yun-Zhu Song and Yi-Syuan Chen and Hong-Han Shuai

National Yang Ming Chiao Tung University, Taiwan

{yunzhusong.eed07g, yschen.eed09g, hhshuai}@nctu.edu.tw

## Abstract

A notable challenge in Multi-Document Summarization (MDS) is the extremely-long length of the input. In this paper, we present an extract-then-abstract Transformer framework to overcome the problem. Specifically, we leverage pre-trained language models to construct a hierarchical extractor for salient sentence selection across documents and an abstractor for rewriting the selected contents as summaries. However, learning such a framework is challenging since the optimal contents for the abstractor are generally unknown. Previous works typically create *pseudo extraction oracle* to enable the supervised learning for both the extractor and the abstractor. Nevertheless, we argue that the performance of such methods could be restricted due to the insufficient information for prediction and inconsistent objectives between training and testing. To this end, we propose a loss weighting mechanism that makes the model aware of the unequal importance for the sentences not in the pseudo extraction oracle, and leverage the fine-tuned abstractor to generate summary references as auxiliary signals for learning the extractor. Moreover, we propose a reinforcement learning method that can efficiently apply to the extractor for harmonizing the optimization between training and testing. Experiment results show that our framework substantially outperforms strong baselines with comparable model sizes and achieves the best results on the Multi-News, Multi-XScience, and WikiCatSum corpora.<sup>1</sup>

## 1 Introduction

Neural multi-document summarization has drawn a lot of attention due to the wide applications, e.g., Wikipedia generation (Liu et al., 2018), news digest (Fabbri et al., 2019), or related-work section generation (Lu et al., 2020). Through con-

catenating the document clusters, it is applicable to adopt the approaches from single-document summarization (Zhang et al., 2020a; Fabbri et al., 2019) under multi-document settings. However, the input length from multiple documents is typically long, which might be computationally infeasible for Transformer-based models (Vaswani et al., 2017). One possible solution is to set a fixed length and truncate the input sequence. However, the truncation leads to information loss and performance drop. Moreover, even if the model can take such a long sequence as input, the attentions could be dispersed over long sequences (Yang et al., 2018) and further degrade the performance (Jin et al., 2020; Liu and Lapata, 2019a; Cohan et al., 2018).

To tackle the long sequence problem, another possible solution is to leverage the extract-then-abstract architecture in single document summarization (Pilault et al., 2020; Chen and Bansal, 2018; Gehrmann et al., 2018). In such an architecture, an extractor first selects salient contents from the documents, and an abstractor is applied to rewrite the selected contents to coherent summaries. However, several issues arise for directly applying the procedure in MDS. 1) *Long-sequence problem for the extractor*. Although the input length for the abstractor can be controlled through the content selection and thus enables the usage of Transformer models (Pilault et al., 2020), it is still inevitable for the extractor to process the complete input documents. Previous methods mainly use LSTM-based extractors (Pilault et al., 2020; Chen and Bansal, 2018; Gehrmann et al., 2018). However, this fails to leverage knowledge of pre-trained language models. 2) *Suboptimal pseudo oracles*. Most summarization corpora do not contain the oracle for extraction. As an alternative, previous works typically generate *pseudo extraction oracle*, or simply called *pseudo oracle*, through a greedy process (Liu and Lapata, 2019b; Chen and Bansal, 2018). Specifically, for each iteration in the process,

<sup>1</sup>The implementation code and trained models are available at <https://github.com/yunzhusong/NAACL2022-REFLECT>

candidate sentences are individually concatenated with previously-selected sentences and compute the ROUGE (Lin, 2004) scores with the human-written summaries. The top-scored sentence is iteratively selected until no sentence could further improve the ROUGE score. In practice, there are different ways for scoring. For examples, Liu and Lapata (2019b) use the average of ROUGE-1 and ROUGE-2 F1 scores while Chen and Bansal (2018) use ROUGE-L recall. The variants of design could cause the extractor to behave differently. In our study, we also found that using the ROUGE precision metric leads to much less extraction than using the recall metric. This implies that the pseudo oracles are suboptimal, and learning the extractor fully relies on the pseudo oracles could restrict the performance. How to alleviate the negative effects of pseudo oracles remains open. 3) *Insufficient information for the extractor*. Even if the pseudo oracles are good enough to train the abstractor well, learning a precise extractor is still challenging. The problem lies in that a pseudo oracle is derived from a specific summary. However, there are potentially multiple valid summaries given the documents. To select salient sentences, the extractor is required to implicitly infer the underlying summary used for oracle construction, which is difficult due to the lack of evidence. 4) *Inconsistent objectives for the extractor*. With pseudo oracles, the extractor is learned to select a set of sentences that has high lexical similarity to the summaries without redundancy. However, the goal of the extractor in test-time is to provide inputs for the abstractor that non-overlapping lexical may still be valuable. In other words, the objective for the extractor in training is inconsistent with the one in testing.

To address these issues, we propose the **RE**ferenced **FLE**xible Extraction with **CrediT**-Awareness (REFLECT) for MDS. For the first problem, we propose a Transformer-based hierarchical extractor that contains the token- and sentence-level feature encoders. Both the encoders are initialized with pre-trained language models to utilize the pretext knowledge. For the second problem, we propose Pseudo Oracle Relaxation (POR) to render the model aware of the unequal importance for the non-oracle sentences. This mechanism encourages the model to emphasize the precision for critical sentences with either high or low lexical similarity to the summaries, and avoids the confusion arising from the different labels for similar

sentences. For the third problem, we propose Summary Referencing (SR) to leverage the fine-tuned abstractor for providing additional learning signals to evidence extraction prediction. The summary reference serves as an approximation for the human-written summary while being able to generalize for testing. For the fourth problem, we propose Credit-Aware Self-Critic (CASC) learning to fine-tune for matching the objective between training and testing. Different from previous methods that assign an identical reward for all actions, we reallocate the rewards based on the impacts of actor explorations.

The contributions are summarized as follows:

- We leverage pre-trained language models to propose an extract-then-abstract framework, which contains a hierarchical extractor that efficiently handles long inputs while utilizing pretext knowledge.
- We investigate the problems for typical learning paradigms of the extractor and propose a framework, named REFLECT, to further improve the extractor performance. The studies on pseudo oracles also provide valuable insights for extract-then-abstract frameworks.
- Experimental results on Multi-News, Multi-XScience, and WikiCatSum corpora demonstrate that REFLECT outperforms the state-of-the-art models with comparable sizes.

## 2 Related Work

Early attempts for MDS focus on extracting sentences through statistical methods (Goldstein et al., 2000; Erkan and Radev, 2004; Wan and Yang, 2006, 2008). For example, Goldstein et al. (2000) extend Maximal Marginal Relevance (MMR) method to select sentences that are relevant to the query and novel across different documents. Erkan and Radev (2004) leverage sentence relations with graph structures that represent pairwise sentence similarities, and apply PageRank (Page et al., 1999) algorithm to extract sentences given the query document. However, extractive methods often suffer the coherence problem (Wu and Hu, 2018). Therefore, instead of directly extracting sentences from the articles, abstractive methods that can rewrite the articles achieve great success with the advantages of large annotated corpora (Pang et al., 2021; Zhou et al., 2021; Liu et al., 2021a; Zhong et al., 2020; Li et al., 2020; Liu and Lapata, 2019a).

Directly operating on the long inputs in MDS often leads to model degradation (Jin et al., 2020). One of the promising solutions is to leverage the extract-then-abstract architectures in single document summarization (Chen and Bansal, 2018; Gehrmann et al., 2018). For example, Gehrmann et al. (2018) propose an LSTM-based word-level extractor to choose phrases from the document, and apply an abstractor with the copy mechanism to generate summaries given the selected contents. Pilault et al. (2020) also apply an LSTM-based sentence extractor to select contents, but further leverage a Transformer decoder to improve the performance.

Although introducing Transformer architectures provides improvements, the input length is typically limited according to the computation and memory overhead. A recent line of studies has been proposed to alleviate such problems (Bražinskas et al., 2021; Beltagy et al., 2020; Liu and Lapata, 2019a; Jin et al., 2020). Liu and Lapata (2019a) propose to rank the candidate input paragraphs, and only concatenate the top few as the inputs for the abstractor. Jin et al. (2020) encode multiple documents in different granularity including token-level, sentence-level, and document-level. Pasunuru et al. (2021) design a graph encoder parallel with standard encoder to provide inter-document information and integrate the pre-trained BART (Lewis et al., 2020) with local attention mechanism (Beltagy et al., 2020) to overcome the long input problem. Bražinskas et al. (2021) apply policy gradient optimization to learn a train-time selector using lexical features pre-computed from source texts and gold summaries as inputs. Due to the lack of gold summaries in testing, a test-time selector is learned using lexical features solely from source texts as inputs and predictions from the train-time selector as learning targets. Comparatively, our methods thoroughly leverage the language models with an extract-and-abstract framework that enjoys the full capacities for both the extractor and the abstractor. The introduced summary referencing further improves the extractor with additional signals in both train- and test-phase.

### 3 Methodology

Utilizing large pre-trained language models brings great benefits for text generation problems. However, the input length of multi-document summarization is typically long, which makes large pre-

trained language models, *e.g.*, Transformer-based models, inefficient for processing. To match the length constraint, document-level truncation (Fabri et al., 2019) has been widely-used. However, the truncation could inevitably cause information loss and degrade the performance. To solve the problem, we propose to leverage the extract-then-abstract framework. However, as described in the introduction, four challenges arise while learning such a framework. We first present our architecture to solve the long sequence problem and elaborate on the proposed methods to tackle the challenges.

#### 3.1 Hierarchical Summarizer

Figure 1 illustrates the proposed REFLECT framework. Specifically, based on the Transformer (Vaswani et al., 2017), our architecture contains a hierarchical extractor (H-EXT) and an abstractor (ABS). The hierarchical extractor is composed of a token-level encoder (T-ENC), a sentence-level encoder (S-ENC), and a sentence selector (SS). Considering a training example  $(x, y)$ ,  $x = \{x_1, x_2, \dots, x_N\}$  is the concatenation of multiple documents that jointly consists of  $N$  complete sentences and  $y$  is the corresponding human-written summary. Let  $M$  denote the allowed input length in Transformer-based models, we split all sentences into  $K$  disjoint sets  $\{h_k\}_{k=1}^K$  such that each set  $h_k$  consists of sentences whose total number of tokens matches the length constraint  $M$ . The hierarchical extractor then predicts a set containing the indices of selected sentences (denoted by  $\hat{e}$ ), which is used to retrieve salient contents from  $x$  for the abstractor to produce summaries. Specifically, the hierarchical extractor can be expressed as follows.

$$\begin{aligned} \hat{e} &= \text{H-EXT}(h) \\ &= \text{SS}(\text{S-ENC}(\bar{\oplus}\{\text{T-ENC}(h_k)\}_{k=1}^K)), \end{aligned} \quad (1)$$

where  $\bar{\oplus}$  is the operation of taking average for hidden states of tokens within sentences and concatenating the averaged results to obtain the sentence-level representations. The hierarchical summarizer (H-SUM) can then be expressed as:

$$\text{H-SUM}(h) = \text{ABS}(\oplus\{w_{x_i} | i \in \hat{e}\}), \quad (2)$$

where  $\oplus$  is the concatenation operation.

#### 3.2 Pseudo Oracle Relaxation

The sentence selection process of the extractor is typically formulated as a supervised classification

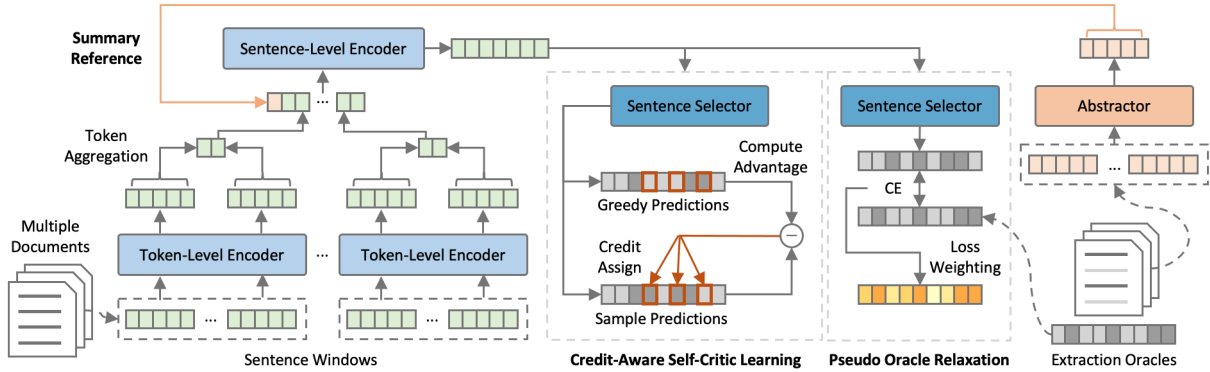


Figure 1: Framework of REFLECT. The illustration of proposed Pseudo Oracle Relaxation (POR), Summary Referencing (SR), and Credit-Aware Self-Critic (CASC) learning are indicated with bold fonts.

problem. However, most summarization corpora do not provide such annotations. As an alternative, it is common to create pseudo oracles through a greedy process with the human-written summaries (Liu and Lapata, 2019b; Chen and Bansal, 2018). However, such methods only provide suboptimal solutions since they are limited by the design of the greedy algorithm. The underlying patterns of true oracles could be too complicated to be designed precisely. The performance of the extractor is thus restricted since typical Maximal Likelihood Estimation (MLE) methods fully depend on pseudo oracles. For example, consider a case that there are only three sentences  $x_a, x_b$  and  $x_c$ , and the ROUGE scores between sentences and summary are ranked as  $x_a > x_b \gg x_c$ . Assume that only  $x_a$  is chosen as the pseudo oracle. If the extractor further selects one additional sentence for improving the results, it is expected that the combination of  $(x_a, x_b)$  could be better than  $(x_a, x_c)$  generally. However, MLE-based methods do not consider such discrepancy according to the pseudo oracle.

Therefore, we propose a loss weighting mechanism to further consider the lexical similarity for the non-oracle sentences during the learning process. Consider an input  $x = \{x_1, x_2, \dots, x_N\}$  and the corresponding pseudo oracles  $e$ , the input sentences can be separated into two sets of pseudo oracles and non-pseudo oracles,  $S = \{x_i | i \in e\}$  and  $S^c = \{x_i | i \in e^c\}$ , respectively. We maintain the loss weights for sentences in  $S$  as one, while modifying the loss weights for the sentences in  $S^c$  as follows:

$$w_i = (1 - \text{ROUGE-1}(x_i, y))^\gamma, \quad (3)$$

where  $\gamma$  is a hyper-parameter controlling the weighting scales. To stabilize the training, we fur-

ther shift loss weights of  $S^c$  such that the maximum value is one. This weighting mechanism emphasizes the predictions of both pseudo oracles and the low-ROUGE sentences, which have more impact on the performance. It relaxes the constraint from the binarized oracles to make the model aware of the differences between candidate sentences. The objective for the hierarchical extractor can then be written as:

$$L = - \sum_{x_i \in x} w_i \log \left( \frac{\exp(z_i^{1_S(i)})}{\exp(z_i^0) + \exp(z_i^1)} \right), \quad (4)$$

where  $z_i^0$  and  $z_i^1$  are the binary logits for the  $i$ -th sentence, and  $1_S(i)$  is the indicator function of oracle sentences  $S$ .

### 3.3 Summary Referencing

The pseudo oracle relaxation described in Section 3.2 makes the extractor focus on the sentences that are more important in terms of the lexical similarity. In other words, it considers the discrepancy between non-oracle sentences instead of treating them equally. However, the ambiguity could still exist between the oracle and non-oracle sentences. For example, consider a case where there are two sentences  $x_a$  and  $x_b$ , and only  $x_a$  is chosen as the oracle during the greedy process. The ROUGE scores between the sentences and the summaries are roughly the same for the two sentences. Learning with such oracles may confuse the model since the positive and negative sentences are similar but with a large loss difference. Therefore, we propose to provide summary references for the extractor to further reason the selection. Since the summary is only available in the training stage, we leverage a fine-tuned abstractor to provide such summary references. With such a mechanism, the operations

of the hierarchical extractor are further revised as:

$$\begin{aligned} r &= \text{ABS}(\oplus\{w_{x_i} | i \in \hat{e}\}), \\ \hat{e} &= \text{H-EXT}(h, r) \\ &= \text{SC}(\text{S-ENC}(\oplus(\{\text{T-ENC}(r)\} + \\ &\quad \{\text{T-ENC}(h_k)\}_{k=1}^K))). \end{aligned} \quad (5)$$

Compared to directly using human-written summaries as references when training, which may cause the extractor highly rely on the reference signal, the generation results from the abstractor could serve as good approximations and provide the generalization from training to testing.

### 3.4 Credit-Aware Self-Critic Learning

With the objective of maximal likelihood estimation, the extractor is learned to select sentences that are jointly have high lexical similarity to the human-written summary. However, in extract-then-abstract frameworks, the required objective of extractor is to provide salient contents that can maximize the generation quality after rewritten by the abstractor. In other words, the objective for the extractor is inconsistent between training and testing. Therefore, we further propose a reinforcement learning method that can directly optimize with the test-time objective to bridge the gaps.

Specifically, we formulate the extraction of sentences as a single-round Combinatorial Multi-Armed Bandit (CMAB) problem (Chen et al., 2016). A general CMAB problem can be modeled as a tuple  $(E, \mathcal{F}, P, R)$ , where  $E = \{1, 2, \dots, N\}$  is a set of  $N$  arms,  $\mathcal{F} \subseteq 2^E$  is a set of subsets of  $E$ ,  $P$  is a probability distribution over  $[0, 1]^N$ , and  $R$  is a reward function defined on  $[0, 1]^N \times \mathcal{F}$ . At the  $t$ -th round, the agent pulls a subset of arms  $S^t \in \mathcal{F}$ , and produce stochastic outcomes  $M^t = (M_1^t, M_2^t, \dots, M_N^t) \sim P$ , where  $M_i^t$  is the outcome of  $i$ -th arm. With a realization of outcomes  $m = (m_1, m_2, \dots, m_N)$ , the agents then receive a reward of  $R(m, S)$ . The goal of the agent is to maximize the expected cumulative reward in  $T$  rounds, which is  $\mathbb{E}_{\{M^t\}_{t=1}^T} [\sum_{t=1}^T R(M^t, S^t)]$ .

For the sentence extraction problem, we consider each sentence as an arm, and pull multiple arms (i.e., select multiple sentences) only for a single round. We solve this problem with the self-critic learning framework (Rennie et al., 2017). Specifically, we consider the sentence-level encoder with the sentence selector that as the agent (jointly parameterized by  $\theta$ ). The agent takes the sentence representations  $x = \{x_1, x_2, \dots, x_N\}$  from the token-

level encoder, and selects a set of sentences  $S$  according to a policy  $\pi_\theta(\cdot)$  as:

$$\begin{aligned} m_i &\sim \text{Bern}\left(\frac{\exp(z_i^1)}{\exp(z_i^0) + \exp(z_i^1)}\right), \\ S &= \pi_\theta(x) = \{i | m_i = 1\}, \end{aligned} \quad (6)$$

where  $z_i^0$  and  $z_i^1$  are the binary logits for the  $i$ -th sentence. To reduce the variance during learning, we introduce a baseline term through another policy  $\tilde{\pi}_\theta$  as:

$$\begin{aligned} \tilde{m}_i &= \frac{\exp(z_i^1)}{\exp(z_i^0) + \exp(z_i^1)}, \\ \tilde{S} &= \tilde{\pi}_\theta(x) = \{i | \tilde{m}_i > 0.5\}, \end{aligned} \quad (7)$$

which is essentially a greedy policy as the taken actions are not explored. The two policies are derived from the agents with the same parameters, and the learning process is thus self-critical. We define the reward as the performance of generation from an abstractor using the selected sentences  $S$  as the input, i.e.,

$$R(S) = \text{ROUEG-L}(\text{ABS}(S), y), \quad (8)$$

where  $S$  is produced from the outcomes  $m$ . The advantage  $a$  can be written as:

$$a = R(S) - R(\tilde{S}). \quad (9)$$

We then optimize the agent through gradient descent with the objective function as:

$$L = - \sum_{i \in E} a \log\left(\frac{\exp(z_i^{1_S(i)})}{\exp(z_i^0) + \exp(z_i^1)}\right), \quad (10)$$

where  $\mathbf{1}_S(i)$  is the indicator function of selected sentence set  $S$ . However, in such a learning objective, the advantages are applied uniformly to update all actions, which could make learning difficult since the sentences number is typically large in our setting. Thus, we propose to specifically credit the advantages to the selections that are distinct between policies  $\pi_\theta$  and  $\tilde{\pi}_\theta$ , and the objective function  $L^{\text{credit}}$  can be written as:

$$\begin{aligned} L^{\text{credit}} &= \\ &- \sum_{i \in E} \mathbf{1}_{S \cap \tilde{S}}(i) a \log\left(\frac{\exp(z_i^{1_S(i)})}{\exp(z_i^0) + \exp(z_i^1)}\right). \end{aligned} \quad (11)$$

Different with previous work (Chen and Bansal, 2018) that formulates the sentence selections as a sequence of decisions and uses LSTM-based agent for learning, our formulation enables the usage of Transformer-based models for a better efficiency.

Model	ROUGE-1	ROUGE-2	ROUGE-L	Average	Average Improvement
Hierarchical-Transformer (Liu and Lapata, 2019a)	42.36	15.27	22.08	26.57	+17.91%
PG-BRNN (Gehrmann et al., 2018)	43.77	15.38	20.84	26.66	+17.52%
Hi-MAP (Fabbri et al., 2019)	44.17	16.05	21.38	27.20	+15.18%
CTF-DPP (Perez-Beltrachini and Lapata, 2021)	45.84	15.94	21.02	27.60	+13.51%
GraphSum (Li et al., 2020)	45.02	16.69	22.50	28.07	+11.61%
GraphSum + RoBERTa (Li et al., 2020)	45.87	17.56	23.39	28.94	+8.26%
Highlight-Transformer (Liu et al., 2021a)	44.62	15.57	-	-	-
MatchSum (Zhong et al., 2020)	46.20	16.51	-	-	-
PEGASUS (Zhang et al., 2020a)	47.52	18.72	<b>24.91</b>	30.38	+3.13%
BART-Long (Pasunuru et al., 2021)	48.54	18.56	23.78	30.29	+3.43%
BART-Long-Graph (Pasunuru et al., 2021)	49.24	18.99	23.97	30.73	+1.95%
REFLECT (MLE)	48.16±0.01	18.87±0.01	23.78±0.17	30.27	+3.50%
REFLECT (CASC)	<b>49.27±0.06</b>	<b>19.96±0.03</b>	24.76±0.09	<b>31.33</b>	-

Table 1: Performance of REFLECT with various baselines on Multi-News corpus. The results of Hierarchical Transformer, HiMAP, and PG-BRNN are copied from Li et al. (2020). The rest baseline results are from the original papers. Best ROUGE scores are bolded. Performance of REFLECT is reported with five runs.

## 4 Experimental Results

### 4.1 Settings

We compare REFLECT with several strong baselines (Liu and Lapata, 2019a; Gehrmann et al., 2018; Fabbri et al., 2019; Perez-Beltrachini and Lapata, 2021; Li et al., 2020; Liu et al., 2021a; Zhong et al., 2020; Zhang et al., 2020a; Pasunuru et al., 2021) on **Multi-News** (Fabbri et al., 2019), **Multi-XScience** (Lu et al., 2020) and **WikiCatSum** (Perez-Beltrachini et al., 2019) corpora, derived from news, academic domains and Wikipedia, respectively. Due to space limit, the results of Multi-XScience and WikiCatSum are provided in the Appendix A. For evaluation, we use ROUGE F1 metrics (Lin, 2004), BERTScore (Zhang et al., 2020b), and factual consistency evaluated with FactCC (Kryscinski et al., 2020) to investigate performance from different perspectives. In our architecture, the hierarchical extractor is initialized by RoBERTa-base (Liu et al., 2020) containing 12 attention layers. We take the first  $l$  layers and the rest layers ( $12-l$ ) as the token- and sentence-level encoder, respectively. The input length limit  $M$  is set to 512. The abstractor is a sequence-to-sequence model initialized by BART (Lewis et al., 2020). To make the CASC computationally efficient, we use the BART-base as the abstractor that provides rewards for the extractor during training, while exploiting the BART-large in the testing. We provide more implementation details in Appendix C.

### 4.2 Main Results

Table 1 summarizes the performance of REFLECT with various baselines on Multi-News

corpus with the rows sorted by the average of ROUGE 1, 2, and L scores (Lin, 2004).<sup>2</sup> The results demonstrate that REFLECT outperforms all baselines. PG-BRNN and Hi-MAP both use an LSTM-based point generator for sentence selection and apply a decoder for the generation. However, the performance is limited since LSTM could still suffer long-term dependency problem (Trinh et al., 2018). Hierarchical-Transformer (row 1) further uses an LSTM-based ranker to select paragraphs through predicting the ROUGE scores with the summaries for each paragraph, and uses a Transformer-based hierarchical encoder to capture the local and global information. However, the estimation of the ROUGE scores could be difficult since there are multiple valid summaries. In contrast, REFLECT further applies pre-trained language models to extract the hierarchical features, and the proposed SR method explicitly provides the extractor with the evidence for selecting more salient contents to be rewritten from the abstractor. Moreover, GraphSum (row 5 & 6) leverages graph structures to explore the relations of paragraphs in the encoder, and Highlight-Transformer (row 7) specifically assigns higher attention weights for key phrases. Both the methods implicitly provide additional selective information for the decoder. REFLECT realizes such a content selection pro-

<sup>2</sup>We do not compare ROUGE-L results with Highlight-Transformer and MatchSum because 1) Highlight-Transformer does not report ROUGE-L results and 2) MatchSum reports summary-level ROUGE-L (ROUGE-LSum) results which are different from sentence-level ROUGE-L used here (Lin, 2004). Nevertheless, ROUGE-LSum results of REFLECT in Table 3 still show the superiority over MatchSum (45.05 v.s. 41.89).

cess in the data-level through the two-stage design. The proposed POR and SR methods alleviate the negative effects from the pseudo oracles, and thus provide complete information for the abstractor to rewrite. The results of PEGASUS (row 9) show the benefits of pre-trained sequence-to-sequence language models, which are also leveraged in REFLECT. The BART-Long-Graph (row 10 & 11) further combines the previous advantages on graph structure and pre-trained language models to achieve strong performance. REFLECT outperforms it in terms of the ROUGE-2 and ROUGE-L scores especially, probably due to the local attention mechanism used in BART-Long-Graph. Although such a design reduces the complexity to accommodate longer sequences, the trade-off for the attention capacity still restricts the model from generating more coherent summaries. Due to the space limit, please refer to Appendix A for results on Multi-XScience and WikiCatSum.

In addition to the evaluation of lexical overlaps, Table 2 shows the performance of semantic and factual consistency, which are important for summarization applications that will influence the public. The results demonstrate that REFLECT can improve factual consistency through the hierarchical architecture, *i.e.*, the architecture enables selection of useful information from multiple documents, while still maintaining the semantics in generations.

Model	BERTScore	Factual Consistency
BART-base	0.870	79.7
CTF-DPP	0.852	81.9
REFLECT	<b>0.871</b>	<b>82.2</b>

Table 2: BERTScore and factual consistency evaluated with FactCC on Multi-News corpus.

### 4.3 Ablations and Analyses

In this subsection, we perform ablations and analyses for REFLECT, and more results can be found in Appendix B.

**Effect of Pseudo Oracle Relaxation (POR).** We start the discussions with the results using MLE objectives. As shown in Table 3, when applying POR to the base method (row 1 & 2), the extraction recall increases while the precision decreases. Pseudo oracles generally contain sentences with high ROUGE scores to the summaries, and POR encourages the selection for such sentences by as-

signing smaller loss weights to increase the recall. From the abstraction performance, it shows that such selection preference could provide more concise information (even though not in pseudo oracles) as the input for improvements. Under the combination of SR (row 3 & 4), the precision is further improved and achieves a higher overall F1 score due to more information for selection. However, the abstraction results show that the performance is actually inferior to the one only using SR, which meets our suggestion that the training objective of the extractor is inconsistent in the test-time. The proposed CASC learning further overcomes such problems and improves the performance by integrating POR and SR (row 6 & 10).

**Effect of Summary Referencing (SR).** We first investigate the MLE results. With SR, the extractor is learned to select sentences given the approximations of summaries generated from the abstractor. The extraction recall is thus enhanced with such references. However, the trade-off between recall and precision still exists due to the discrepancy between ground-truth summaries and generation results (row 1 & 3). Although ground-truth can be used as references to increase extraction performance during training, it fails to generalize in testing due to the distributional difference between them. The abstraction performance also shows SR makes significant improvements (row 3) and also benefits from the combination with CASC (row 6 & 8).

**Effect of Credit-Aware Self-Critic (CASC) Learning.** Consider the case that applying CASC with POR and SR respectively (row 5 & row 6, row 7 & row 8), the results demonstrate that CASC substantially outperforms the Self-Critic (SC) learning methods, even when combining the usage of POR and SR (row 9 & row 10). We could further investigate the difference between MLE and RL methods through the visualization in Figure 2. It shows that, although having lower extraction performance, the RL methods can improve over the MLE methods in final abstraction performance due to the consistent objective between training and testing. CASC further improves over SC through explicitly assigning advantage to the actions that have the credits for exploration. Such design is critical for the long input sequences in multi-document summarization and can potentially be applied to other applications.

MLE		RL		Extraction Performance			Abstraction Performance				
SR	POR	SR	CA	Precision	Recall	F1	R-1	R-2	R-L	R-LSum	Average
		-	-	0.6489	0.4218	0.5113	47.60	18.48	23.89	43.22	33.29
	✓	-	-	0.4618	0.7684	0.5769	47.92	18.76	23.57	43.87	33.53
✓		-	-	0.4755	0.6275	0.5410	<b>48.40</b>	<b>18.95</b>	<b>24.08</b>	44.03	<b>33.86</b>
✓	✓	-	-	0.4926	0.7112	0.5820	48.16	18.87	23.61	<b>44.06</b>	33.68
	✓			0.6161	0.4704	0.5335	48.59	19.04	24.20	44.21	34.01
	✓		✓	0.5344	0.5948	0.5630	48.91	19.77	24.46	44.80	34.49
✓		✓		0.5933	0.5328	0.5614	48.83	19.50	24.36	44.57	34.32
✓		✓	✓	0.5442	0.5942	0.5681	49.04	19.80	24.50	44.94	34.57
✓	✓	✓		0.6007	0.5170	0.5557	48.84	19.41	24.34	44.56	34.29
✓	✓	✓	✓	0.5873	0.5108	0.5464	<b>49.27</b>	<b>19.96</b>	<b>24.76</b>	<b>45.04</b>	<b>34.76</b>

Table 3: Ablation of REFLECT with Pseudo Oracle Relaxation (POR), Summary Referencing (SR), and Credit-Aware Self-Critic (CASC) learning. ROUGE score is abbreviated as  $R$ .

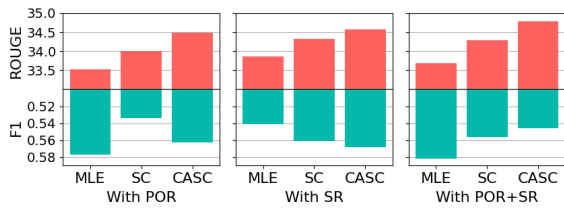


Figure 2: Comparisons of MLE and RL results. For clarity, the displayed range for average ROUGE scores and extraction F1 scores are limited in  $[33,35]$  and  $[0.5,0.59]$ , respectively.

#### 4.4 Effects of Summary Reference

REFLECT uses the BART-large fine-tuned with truncated articles (BART-large-A) to provide summary references for both training and testing. Here, we study the effect of applying different summary referencing strategies in test-time. To investigate such effect across models, we leverage PEGASUS (Zhang et al., 2020a) and fine-tune it with Multi-News corpus to produce summary references. Furthermore, we also conduct experiments that use the BART-large fine-tuned with pseudo oracle sentences (BART-large-O) and the ground-truth summaries. *Note that these two methods are prohibited in practical usage since they require annotations on testing data.* The results in Table 4 manifest that using references from another model that is different from the training one does not significantly affect the performance. For the case using oracle inputs, the performance is similar to the ones directly using articles. This suggests that our method is stable for the practical scenario. The result of the ground-truth also shows that the reference quality could still affect the performance, which is left as a future research direction.

Reference Source	R-1	R-2	R-L	Average
PEGASUS	49.28	19.99	24.84	34.78
BART-large-A	49.29	19.99	24.83	34.79
BART-large-O	49.29	19.98	24.81	34.78
Ground-Truth	49.38	20.05	24.88	34.85

Table 4: Performance of different summary referencing strategies in REFLECT. We also present the results of BART models trained by truncated articles (A) and pseudo oracles (O). ROUGE score is abbreviated as  $R$ .

## 5 Conclusion

In this work, we present an effective extract-then-abstract framework, named REFLECT, for MDS. We utilize large pre-trained language models to construct a hierarchical extractor for solving the long sequence problem. Moreover, we investigate current learning paradigms for such frameworks and find that entirely relying on the pseudo oracles produced via a greedy process could hinder the performance. Therefore, we propose three corresponding techniques (POR, SR, and CASC) to overcome the issues. The experimental results not only show that REFLECT outperforms the state-of-the-art models on Multi-News, Multi-XScience, and WikiCatSum corpora, but also demonstrate that bridging the gap between training and testing is significant. Also, we present extensive studies to motivate more investigations. Finally, we consider further exploring the interactions between the extractor and abstractor, including iteratively providing a better reference or reusing the extraction predictions as the training signals for the abstractor, and study how to build a more efficient reference for providing the extraction evidence.



## Acknowledgements

This work is supported in part by the Ministry of Science and Technology (MOST) of Taiwan under the grants MOST-109-2221-E-009-114-MY3, MOST-110-2221-E-001-001 and MOST-110-2221-E-A49-164. This work was also supported by the Higher Education Sprout Project of the National Yang Ming Chiao Tung University and Ministry of Education (MOE), Taiwan. We are grateful to the National Center for High-performance Computing for computer time and facilities.

## References

- Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.
- Arthur Braźinskas, Mirella Lapata, and Ivan Titov. 2021. Learning opinion summarizers by selecting informative reviews. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9424–9442, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Wei Chen, Wei Hu, Fu Li, Jian Li, Yu Liu, and Pinyan Lu. 2016. Combinatorial multi-armed bandit with general reward functions. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.
- Yen-Chun Chen and Mohit Bansal. 2018. Fast abstractive summarization with reinforce-selected sentence rewriting. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 675–686, Melbourne, Australia. Association for Computational Linguistics.
- Arman Cohan, Franck Dernoncourt, Doo Soon Kim, Trung Bui, Seokhwan Kim, Walter Chang, and Nazli Goharian. 2018. A discourse-aware attention model for abstractive summarization of long documents. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 615–621, New Orleans, Louisiana. Association for Computational Linguistics.
- Günes Erkan and Dragomir R. Radev. 2004. Lexrank: Graph-based lexical centrality as salience in text summarization. *J. Artif. Int. Res.*, 22(1):457–479.
- Alexander Fabbri, Irene Li, Tianwei She, Suyi Li, and Dragomir Radev. 2019. Multi-news: A large-scale multi-document summarization dataset and abstractive hierarchical model. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1074–1084, Florence, Italy. Association for Computational Linguistics.
- Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N. Dauphin. 2017. Convolutional sequence to sequence learning. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1243–1252. PMLR.
- Sebastian Gehrmann, Yuntian Deng, and Alexander Rush. 2018. Bottom-up abstractive summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4098–4109, Brussels, Belgium. Association for Computational Linguistics.
- Jade Goldstein, Vibhu Mittal, Jaime Carbonell, and Mark Kantrowitz. 2000. Multi-document summarization by sentence extraction. In *Proceedings of the 2000 NAACL-ANLP Workshop on Automatic Summarization - Volume 4*, NAACL-ANLP-AutoSum '00, page 40–48, USA. Association for Computational Linguistics.
- Hanqi Jin, Tianming Wang, and Xiaojun Wan. 2020. Multi-granularity interaction network for extractive and abstractive multi-document summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6244–6254, Online. Association for Computational Linguistics.
- Wojciech Kryscinski, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. Evaluating the factual consistency of abstractive text summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9332–9346, Online. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Wei Li, Xinyan Xiao, Jiachen Liu, Hua Wu, Haifeng Wang, and Junping Du. 2020. Leveraging graph to improve abstractive multi-document summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6232–6243, Online. Association for Computational Linguistics.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Peter J. Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. 2018. Generating wikipedia by summarizing long sequences. In *International Conference on Learning Representations*.

- Shuaiqi Liu, Jiannong Cao, Ruosong Yang, and Zhiyuan Wen. 2021a. [Highlight-transformer: Leveraging key phrase aware attention to improve abstractive multi-document summarization](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 5021–5027, Online. Association for Computational Linguistics.
- Yang Liu and Mirella Lapata. 2019a. [Hierarchical transformers for multi-document summarization](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5070–5081, Florence, Italy. Association for Computational Linguistics.
- Yang Liu and Mirella Lapata. 2019b. [Text summarization with pretrained encoders](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3730–3740, Hong Kong, China. Association for Computational Linguistics.
- Yang Liu, Sheng Shen, and Mirella Lapata. 2021b. [Noisy self-knowledge distillation for text summarization](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 692–703, Online. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [RoBERTa: A robustly optimized BERT pretraining approach](#).
- Yao Lu, Yue Dong, and Laurent Charlin. 2020. [MultiXScience: A large-scale dataset for extreme multi-document summarization of scientific articles](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8068–8074, Online. Association for Computational Linguistics.
- Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. 1999. [The pagerank citation ranking: Bringing order to the web](#). Technical Report 1999-66, Stanford InfoLab.
- Richard Yuanzhe Pang, Adam Lelkes, Vinh Tran, and Cong Yu. 2021. [AgreeSum: Agreement-oriented multi-document summarization](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3377–3391, Online. Association for Computational Linguistics.
- Ramakanth Pasunuru, Mengwen Liu, Mohit Bansal, Sujith Ravi, and Markus Dreyer. 2021. [Efficiently summarizing text and graph encodings of multi-document clusters](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4768–4779, Online. Association for Computational Linguistics.
- Laura Perez-Beltrachini and Mirella Lapata. 2021. [Multi-document summarization with determinantal point process attention](#). *J. Artif. Int. Res.*, 71:371–399.
- Laura Perez-Beltrachini, Yang Liu, and Mirella Lapata. 2019. [Generating summaries with topic templates and structured convolutional decoders](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5107–5116, Florence, Italy. Association for Computational Linguistics.
- Jonathan Pilault, Raymond Li, Sandeep Subramanian, and Chris Pal. 2020. [On extractive and abstractive neural document summarization with transformer language models](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9308–9319, Online. Association for Computational Linguistics.
- Steven J. Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. 2017. [Self-critical sequence training for image captioning](#). In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1179–1195, Los Alamitos, CA, USA. IEEE Computer Society.
- Trieu Trinh, Andrew Dai, Thang Luong, and Quoc Le. 2018. [Learning longer-term dependencies in RNNs with auxiliary losses](#). In *Proceedings of the 35th International Conference on Machine Learning (PMLR)*, volume 80, pages 4965–4974.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Xiaojun Wan and Jianwu Yang. 2006. [Improved affinity graph based multi-document summarization](#). In *Proceedings of the Human Language Technology Conference of the NAACL, Companion Volume: Short Papers*, pages 181–184, New York City, USA. Association for Computational Linguistics.
- Xiaojun Wan and Jianwu Yang. 2008. [Multi-document summarization using cluster-based link analysis](#). In *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '08*, page 299–306, New York, NY, USA. Association for Computing Machinery.
- Yuxiang Wu and Baotian Hu. 2018. [Learning to extract coherent summary via deep reinforcement learning](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1).
- Baosong Yang, Zhaopeng Tu, Derek F. Wong, Fandong Meng, Lidia S. Chao, and Tong Zhang. 2018. [Modeling localness for self-attention networks](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4449–4458, Brussels, Belgium. Association for Computational Linguistics.

Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter Liu. 2020a. [PEGASUS: Pre-training with extracted gap-sentences for abstractive summarization](#). In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 11328–11339. PMLR.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020b. [Bertscore: Evaluating text generation with bert](#). In *International Conference on Learning Representations*.

Ming Zhong, Pengfei Liu, Yiran Chen, Danqing Wang, Xipeng Qiu, and Xuanjing Huang. 2020. [Extractive summarization as text matching](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6197–6208, Online. Association for Computational Linguistics.

Hao Zhou, Weidong Ren, Gongshen Liu, Bo Su, and Wei Lu. 2021. [Entity-aware abstractive multi-document summarization](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 351–362, Online. Association for Computational Linguistics.

Fangwei Zhu, Shangqing Tu, Jiaxin Shi, Juanzi Li, Lei Hou, and Tong Cui. 2021. [TWAG: A topic-guided Wikipedia abstract generator](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, pages 4623–4635, Online. Association for Computational Linguistics.

## A Results on Other Corpora

### A.1 REFLECT on Multi-XScience Corpus

In this subsection, we additionally experiment REFLECT with Multi-XScience (Lu et al., 2020) corpus, which is a multi-document summarization corpus for scientific articles. The source of this corpus is a tuple of a paper abstract and abstracts of the citation papers, while the generation target is the corresponding related-work paragraph. The performances of baselines are from Lu et al. (2020)<sup>3</sup>. Table 5 shows that REFLECT achieves a better performance comparing to all baselines, and makes a substantial improvement especially in ROUGE-2 score. We also report the results of BART-large which serves as our abstractor to summarize from the extracted sentences. The improvement shows the benefit of the extraction process after the Credit-Aware Self-Critic Learning.

<sup>3</sup>We can not tell the version of score metric ROUGE-L used in Multi-XScience, and the evaluation code is also unavailable. Therefore, we do not report the ROUGE-L score from Lu et al. (2020)

Model	R-1	R-2	R-L	R-LSum
LEAD	27.46	4.57	-	-
LEXRANK	30.19	5.53	-	-
TEXTRANK	31.51	5.83	-	-
Hi-MAP	31.66	5.91	-	-
PG	34.11	6.76	-	-
BART-large	33.29	8.07	17.31	29.04
REFLECT (MLE)	33.87	8.13	17.20	29.69
REFLECT (CASC)	<b>34.18</b>	<b>8.20</b>	<b>17.42</b>	<b>29.73</b>

Table 5: Performance of REFLECT and baselines on Multi-XScience corpus. We also present the results of BART models trained by truncated articles. The results of baselines are from Lu et al. (2020), where PG represents Pointer-Generator and the Hi-MAP represents the model proposed by Fabbri et al. (2019). ROUGE score is abbreviated as *R*.

### A.2 REFLECT on WikiCatSum Corpus

To further analyze the capability of our proposed model, we also apply REFLECT on a domain-specific corpus, WikiCatSum (Perez-Beltrachini et al., 2019) as shown in Table 6. WikiCatSum is a multi-document summarization dataset derived from WikiSum (Liu et al., 2018) and represents three distinct domains (Animal, Company, and Film). We compare our model with several baselines. TF-S2S is a Transformer sequence-to-sequence model of Liu et al. (2018) and CV-S2D+T is a variant of CV-S2S (Gehring et al., 2017) with a single sequence encoder and a structure decoder. Both CV-S2D+T (Perez-Beltrachini et al., 2019) and TWAG (Zhu et al., 2021) introduce the topic detection model to guide the generation. The decoder of CV-S2D+T is trained to predict the topics as an auxiliary task, while TWAG uses the topic information to group the input paragraphs and encodes them separately. Liu et al. (2021b) apply knowledge distillation to alleviate the problem of single reference in maximum likelihood training, while our model leverages reinforcement learning to the train-test mismatch issue. The results of BART-large models are fine-tuned on WikiCatSum<sup>4</sup>. REFLECT outperforms all baselines in all three domains, especially in the ROUGE-1 score, showing that our generated summaries carry more information. The results between REFLECT (MLE) and REFLECT (CASC) also manifest the effectiveness of Credit-Aware Self-Critic learning for bridging the gap between training and testing.

<sup>4</sup>The models of Liu et al. (2021b) and REFLECT are with pretraining. We exploit the pretrained model from <https://huggingface.co/facebook/bart-large-cnn>.

Animal				
Model	R-1	R-2	R-L	R-LSum
TF-S2S	44.0	28.8	40.0	-
CV-S2D+T	42.7	27.9	37.9	-
TWAG	43.1	24.4	40.9	-
Liu et al. (2021b)	45.9	<b>32.2</b>	41.4	-
<hr/>				
BART-large	46.3	29.2	39.6	44.1
REFLECT (MLE)	46.5	27.1	38.2	43.6
REFLECT (CASC)	<b>48.6</b>	30.2	<b>41.5</b>	<b>46.0</b>
<hr/>				
Company				
Model	R-1	R-2	R-L	R-LSum
TF-S2S	26.0	9.5	20.4	-
CV-S2D+T	27.5	10.6	21.4	-
TWAG	34.1	11.9	<b>31.6</b>	-
Liu et al. (2021b)	33.5	15.0	25.9	-
<hr/>				
BART-large	36.8	15.1	25.6	33.6
REFLECT (MLE)	40.3	15.5	27.0	36.2
REFLECT (CASC)	<b>40.8</b>	<b>15.8</b>	27.5	<b>36.6</b>
<hr/>				
Film				
Model	R-1	R-2	R-L	R-LSum
TF-S2S	36.5	18.8	31.0	-
CV-S2D+T	38.0	21.2	32.3	-
TWAG	40.8	21.2	34.3	-
Liu et al. (2021b)	42.7	26.1	36.8	-
<hr/>				
BART-large	44.3	25.5	35.9	41.7
REFLECT (MLE)	46.7	25.6	36.5	43.2
REFLECT (CASC)	<b>47.6</b>	<b>26.8</b>	<b>37.9</b>	<b>44.1</b>

Table 6: Performance of REFLECT with various baselines on WikiCatSum corpus. The results of Transformer sequence-to-sequence (TF-S2S) and CV-S2D+T are referenced from Perez-Beltrachini et al. (2019), and TWAG is from Zhu et al. (2021). REFLECT outperforms all baseline models in a large margin for all domains, especially for Company and Film in the ROUGE-1 score.

## B Additional Ablations & Analyses

### B.1 Effects of Hierarchical Architecture

In this subsection, we study the effect of hierarchical architecture in the extractor. As described in Section 4.1, the token- and sentence-level encoder jointly contain 12 layers in our settings. Thus, we experiment with different numbers of layers distributed for the token-level encoder. Also, we consider a case for no (0) layer, where we directly aggregate the word embeddings from the token-level encoder as sentence features. Table 7 demonstrates that hierarchical architecture improves the resulted abstraction performance through better sentence representations. The performance is enhanced with the increasing number of layers, while more than 3 layers only makes marginal improvements. Thus, we use 3 token-level layers in this paper.

Layer Number	R-1	R-2	R-L	R-LSum	Average
0	47.94	18.68	23.46	43.87	33.49
1	48.01	18.63	23.53	43.95	33.53
2	48.11	<b>18.87</b>	<b>23.62</b>	44.03	33.66
3	<b>48.16</b>	<b>18.87</b>	23.61	<b>44.06</b>	<b>33.68</b>

Table 7: Performance of different extractor configurations with MLE learning. Note that all configurations share the same amount of learnable parameters. ROUGE score is abbreviated as  $R$ .

### B.2 The Choice of Summary Reference at Training Stage

In this subsection, we study the effect of the summary reference at training stage. As mention in the subsection 3.3, directly using human-written summaries as reference may cause the severe train-test mismatch. To verify the idea, we take the human-written summaries as the summary referencing for the extractor during training, and perform the standard test process, that is utilizing generated summaries as the references. The results shown in Table 8 verify our assumption that training with such explicit signals reduces the ability of the model to generalize. Therefore, in our experiments, we take the generation results of BART-large as the summary referencing of the extractor.

MLE	R-1	R-2	R-L	R-LSum	Average
Ground-truth	47.54	18.60	23.40	43.63	33.29
Generated	48.16	18.87	23.61	44.06	33.68

Table 8: Performance of the MLE training results when taking the ground-truth summary as the summary reference (SR) during training. ROUGE score is abbreviated as  $R$ .

### B.3 Effect of Input Settings for Abstractor Fine-tuning

In REFLECT, we decouple the learning for the extractor and the abstractor. Therefore, we study the effect of the input settings for the fine-tuning of the abstractor. Table 9 shows the performance of BART under different model configurations with either using truncated articles or pseudo oracle sentences as inputs. We find that using pseudo oracles for BART-large fine-tuning can effectively improve the performance even when the testing input is truncated articles. In addition, the results of using pseudo oracles for both training and testing demonstrate the upper bound performance of the extractor

under MLE training, which suggests a promising development of our framework as the extractor can make more precise predictions.

BART	Train/Test	R-1	R-2	R-L	R-LSum
base	A/A	45.71	17.12	23.82	41.54
base	O/A	45.62	16.58	22.99	41.59
base	A/O	49.53	22.02	26.78	45.01
base	O/O	51.93	23.89	27.78	47.42
large	A/A	46.80	18.01	23.80	42.57
large	O/A	47.79	18.37	23.78	43.57
large	A/O	51.01	22.77	27.00	46.53
large	O/O	52.98	24.28	28.02	48.34

Table 9: Performance of BART with different configurations and train/test input settings. The input settings include truncated articles (A) and concatenated sentences with pseudo oracles (O). ROUGE score is abbreviated as  $R$ .

#### B.4 Generation Examples

We demonstrate two examples of generation in Table 10 and Table 11. The results demonstrate that generations from REFLECT-(CASC) could provide more faithful information from multiple documents, which mainly resulted from the better sentence extraction strategy learned through CASC.

### C Implementation Details

All of our experiments are conducted on a single NVIDIA Tesla V100 32GB GPU with PyTorch. The hierarchical extractor is initialized by RoBERTa-base<sup>5</sup>, and the first three layers are exploited as token-level encoder and the rest layers are sentence-level encoder. The BART-base<sup>6</sup> is used as the abstractor to provide the rewards during extractor training, while BART-large<sup>7</sup> is used for generating the final results. We generate the initial summary reference (SR) by the BART-large model. The hyper-parameter  $\gamma$  in POR is set to 10. We use Adam with constant learning  $3e-5$  for optimization, and select the model with highest ROUGE-1 F1 score on validation set. For the pseudo extraction oracle, we greedily select at least 30 sentences from article as the pseudo oracle for the extractor during MLE learning. The selection criteria is based on the average of ROUGE-1 recall and ROUGE-2 recall.

For evaluations, we report ROUGE, BERTScore and factual consistency derived from FactCC frame-

work. For ROUGE, we use *rouge\_score* package<sup>8</sup> to report ROUGE-1, 2, L, and LSum scores. For BERTScore, we use official implementation<sup>9</sup> and report the F1 scores. For factual consistency, we use official implementation of FactCC<sup>10</sup> and follow previous work (Perez-Beltrachini and Lapata, 2021) to calculate the factual consistency for multi-document summarization.

<sup>5</sup><https://huggingface.co/deepset/roberta-base-squad2>

<sup>6</sup><https://huggingface.co/facebook/bart-base>

<sup>7</sup><https://huggingface.co/facebook/bart-large>

<sup>8</sup><https://pypi.org/project/rouge-score>

<sup>9</sup>[https://github.com/Tiiiger/bert\\_score](https://github.com/Tiiiger/bert_score)

<sup>10</sup><https://github.com/salesforce/factCC>

---

**Article:** Wesley Snipes Begins Serving 3-Year Prison Stint. This Wesley Snipes began serving a three-year sentence at a federal prison in Pennsylvania on Thursday for failure to file income tax returns. Snipes, 48, arrived shortly before noon at the Federal Correctional Institution McKean in the tiny northwestern Pennsylvania town of Lewis Run, federal prisons spokesman Ed Ross said. He had been ordered to surrender by noon. The minimum security prison camp is worlds away from the harsh prison fortresses depicted in the Snipes' films 'Undisputed' and 'Brooklyn's Finest.' The minimum-security camp doesn't have fences around its perimeter. The 300 nonviolent inmates live in barracks that feature two-man rooms, daily showers and double-feature movie showings Friday through Sunday. ...

---

**Gold Summary:** – Wesley Snipes enters prison in Pennsylvania today to serve a three-year sentence for tax evasion—and though the minimum-security McKean prison camp isn't exactly a five-star resort, it's not as bad as it could be. The AP gives details on what his life will be like: No fences: But he will submit to five daily head counts, three of which are overnight. Living quarters: He'll have a two-man room in the barracks, where the other 300 nonviolent inmates live. Daily schedule: Wake-up time is 6:35am, and jobs are performed for seven hours per day. Conjugal visits: He'll have to limit himself to just a kiss in the visitors room. Money: He can earn pennies an hour by doing laundry or other chores, and is allowed to spend \$290 a month at the commissary. Entertainment: Double-feature movies are shown Friday through Sunday, but no R-, NC-17-, or X-rated films are screened. Exercise: Sand volleyball, indoor basketball, exercise machines, badminton, bocci, and bridge will be available to him. Medical facilities: The copay at the infirmary is only \$2. Church: Almost any group you can think of, including Wiccans, hold weekly meetings. Click for more, including why Snipes failed to pay taxes—or watch his Tuesday night appearance on Larry King Live.

---

**REFLECT (MLE):** – Wesley Snipes has begun his **three-year prison sentence** for failure to file income **tax** returns. The Blade: Trinity **star** arrived at the Federal Correctional Institution **McKean** in Lewis Run in **Pennsylvania today**, reports **the AP**. The **minimum security prison camp** is worlds away from the harsh prison fortresses depicted in Snipes' **films** Undisputed and Brooklyn's Finest, notes the AP, but the 48-year-old actor will **have to** adjust to earning just **pennies an hour** handling kitchen, **laundry**, and other campus **chores**. "We recognize that he is high profile, but we treat all our **inmates** the same," says a spokeswoman. Snipes, a dues-paying member of a tax-protest **group** that challenges the government's right to collect **taxes**, has been free on bail since his 2008 conviction, but he made a last-minute request for a new trial this week, arguing that the judge erred by not allowing defense attorneys to interview jurors about misconduct allegations. He also told **Larry King** last night that he wasn't nervous about reporting to prison. "Don't send me up the river, yet," he said.

---

**REFLECT (CASC):** – Wesley Snipes has begun serving a **three-year sentence** in a **minimum-security prison in Pennsylvania** for failure to file income **tax** returns. The Blade: Trinity **star** arrived at the Federal Correctional Institution **McKean** yesterday, just before noon, reports **the AP**. He had been ordered to surrender by noon. The minimum security **prison camp** is worlds away from the harsh prison fortresses depicted in the Snipes' **films** Undisputed and Brooklyn's Finest. The **300 nonviolent inmates live in barracks** that **feature two-man rooms, daily showers, and double-feature movie showings Friday through Sunday**. The martial-arts enthusiast can get his **exercise** playing **sand volleyball** or **indoor basketball**, or work out on an elliptical machine or stair climber. Alas, no NC-17, R or X ratings **allowed**, which knocks out much of Snipes' action-heavy repertoire. The most jarring aspect of the celebrity's stay might be the **five daily head counts, three** during the **overnight hours**. And Snipes, who earned a reported \$13 million for the Blade: Trinity sequel, will **have to** adjust to earning just **pennies an hour** handling kitchen, **laundry, or** other campus **chores**. **He can spend just \$290 a month at the prison commissary**. "We recognize that he is high profile, but we treat all our inmates the same," says a spokeswoman. Snipes made a last-minute request for a new trial on Wednesday, arguing that the judge erred by not allowing defense attorneys to interview jurors about misconduct allegations.

---

Table 10: A generation example for REFLECT on Multi-News corpus. Key overlaps between generations and the gold summary are bolded.

---

**Documents:** Vacant lot at 53 New York Avenue NE in Washington, where a government permit has been granted for work connected with Elon Musk's Hyperloop project. (Michael Laris/TWP) It's not much now, just a parking lot with a discarded gin bottle and an old exterminator receipt. But the slice of pavement near the Bureau of Alcohol, Tobacco, Firearms and Explosives in the District could be the gritty precursor to a tunnel network that could propel pods filled with people and speeding platforms topped with Teslas and Toyotas between the nation's capital and New York in 29 minutes. Or it could be just be a parking lot littered with dashed transportation dreams. Electric-car pioneer and space entrepreneur Elon Musk has been touting his vision for a high-speed transportation system since his tweeted announcement last summer that he had received "verbal govt approval" for his tunnel-digging firm, the Boring Company, to build a "NY-Phil-Balt-DC Hyperloop." The Boring Company team has received an early, and vague, building permit from the D.C. government...

---

**Gold Summary:** -- Elon Musk's dream of building a hyperloop that can move people between Washington, DC, and New York City in 29 minutes may be a small step closer to becoming a distant reality. A Nov. 29 permit issued by DC's Department of Transportation allows Musk's Boring Company to dig at an abandoned lot beside a McDonald's in northeast Washington that the company says might eventually become a station on a hyperloop connecting NYC, Philadelphia, Baltimore, and Washington. The system would speed up travel times, using electricity and electromagnetic levitation to send people and cars whipping through tunnels in pods. But before you ditch your car, the Washington Post points out the permit is only for "some preparatory and excavation work" at the site, and more permits are needed before any real construction begins. DC Mayor Muriel Bowser has visited a test tunnel Musk is constructing under Los Angeles. But though "we're open to the concept of moving people around the region more efficiently ... we're just beginning, in the mayor's office, our conversation to get an understanding of what the general vision is for Hyperloop," says Bowser's chief of staff, John Falcicchio. Still, the permit "is comparable to an athlete stretching before a race" and shows Musk is serious about moving forward with his plans for high-speed travel since claiming "verbal" government support for an East Coast hyperloop last July, reports TechCrunch. He's also received a conditional permit to build a tunnel beneath Maryland Route 295 in Hawthorne, per the Post, with Maryland officials saying a leg linking Washington and Baltimore will come first. (Read about the hyperloop pod's first test.)

---

**REFLECT (MLE):** – **Elon Musk's Boring Company has received a permit from the DC government to begin work on a tunneling project in the nation's capital, the Washington Post reports.** The permit, **issued Nov. 29,** is a "very small, very vague **step forward**" for Musk's plan **to build a "NY-Phil-Balt-DC Hyperloop"** tunnel network that would run from **New York to Philadelphia to Baltimore to DC in 29 minutes,** according to the New York Times. But it could also be a "gravel precursor to a **tunnel network**" that could propel **pods** filled with **people and** speeding platforms topped with Teslas and Toyotas between the nation's capital **and New York** in just 29 minutes. Musk's **company** has received backing from the White House **Office** of American Innovation, led by Jared Kushner, as well as **Maryland Gov. Martin O'Malley and DC Mayor Muriel Bowser,** who **visited a tunnel** in California earlier this month to learn more **about the technology the company says** will make tunneling faster and cheaper. The DC **Department of Transportation** is figuring out what other **permits** the company will need to cut under **city roads** and other public spaces. Asked whether the Bowser administration supports the project, he was somewhat upbeat but noncommittal, adding: **"We're open to the concept of moving people around the region more efficiently."**

---

**REFLECT (CASC):** – **Elon Musk's plan to build a high-speed transportation system** that would run from **New York to Philadelphia to Baltimore to DC in 29 minutes** has gotten a "very small, very vague **step forward,**" **the Washington Post reports.** Musk's **Boring Company has received a permit** for preparation and preliminary **excavation of a site in Washington, DC,** next to a **McDonald's** and near the Bureau of Alcohol, Tobacco, Firearms, **and Explosives.** It's not much now, just a parking **lot** with a discarded gin bottle and an old exterminator receipt, but it could be "the gritty precursor to a **tunnel network** that could propel **pods** filled with **people and** speeding platforms topped with Teslas and Toyotas between the nation's capital **and New York** in 29 minutes," **per the Post.** Or it could "be just be a Parking **Lot** littered with dashed transportation dreams." Musk has been touting his **vision for a Hyperloop** system since last summer, when he tweeted that he had received "**verbal government approval**" to build the "NY-Phil-Balt-DC Hyperloop." **Maryland officials** told the Post **that the tunnel would run under Maryland Route 295, with the DC-Baltimore leg being built first.** **DC Mayor Muriel E. Bowser visited the company** in California this month, walking in a **tunnel** to learn more **about the technology the company says** will make tunneling faster and cheaper. The District's **Department of Transportation** is figuring out what other **permits** the company would need to cut under **city roads** and other public spaces, according to **Bowser's chief of staff.** **"We're just beginning, in the mayor's office, our conversation to get an understanding of what the general vision is for Hyperloop,"** he says. Asked whether the Bowser administration supports the project, he says, **"We're open to the concept of moving people around the region more efficiently."**

---

Table 11: A generation example for REFLECT on Multi-News corpus. Key overlaps between generations and the gold summary are bolded.