

# Reward Learning for Efficient Reinforcement Learning in Extractive Document Summarisation

Yang Gao<sup>1\*</sup>, Christian Meyer<sup>2</sup>, Mohsen Mesgar<sup>2</sup> and Iryna Gurevych<sup>2</sup>

<sup>1</sup>Dept. of Computer Science, Royal Holloway, University of London

<sup>2</sup>Ubiquitous Knowledge Processing Lab (UKP-TUDA), Technische Universität Darmstadt

yang.gao@rhul.ac.uk, {meyer,mesgar,gurevych}@ukp.informatik.tu-darmstadt.de

## Abstract

*Document summarisation* can be formulated as a sequential decision-making problem, which can be solved by *Reinforcement Learning (RL)* algorithms. The predominant RL paradigm for summarisation learns a *cross-input* policy, which requires considerable time, data and parameter tuning due to the huge search spaces and the delayed rewards. Learning *input-specific* RL policies is a more efficient alternative but so far depends on handcrafted rewards, which are difficult to design and yield poor performance. We propose RELIS, a novel RL paradigm that learns a reward function with *Learning-to-Rank (L2R)* algorithms at training time and uses this reward function to train an input-specific RL policy at test time. We prove that RELIS guarantees to generate near-optimal summaries with appropriate L2R and RL algorithms. Empirically, we evaluate our approach on extractive multi-document summarisation. We show that RELIS reduces the training time by two orders of magnitude compared to the state-of-the-art models while performing on par with them.

## 1 Introduction

*Extractive document summarization*, as a challenging instance of *natural language generation (NLG)*, is a popular summarisation paradigm, which builds summaries by selecting an appropriate sequence of important phrases or sentences from the input document(s). Extractive summarisation can be formulated as a sequential decision-making problem, and hence can be tackled by the *Reinforcement Learning (RL)* algorithms. RL searches for the (near-)optimal *trajectories* (i.e. sequences of decisions) by directly optimising the objective functions, e.g. the ROUGE metrics [Lin, 2004]. Such objectives are non-differentiable and therefore difficult to be directly optimised by deep neural networks. In addition, RL alleviates the *exposure bias* problem faced by sequential supervised learning paradigms in NLG. Combining RL algorithms such as REINFORCE [Williams, 1992] with neural techniques (e.g. sequence-to-sequence) yields state-of-the-art

performance in summarisation [Narayan *et al.*, 2018; Yao *et al.*, 2018].

Existing RL-based summarisation systems fall into two categories: *cross-input RL* and *input-specific RL*. For cross-input RL (upper part in Fig. 1), at training time, RL agents interact with a *ground-truth reward oracle* for multiple *episodes* so as to learn a *policy* that maximises the accumulated rewards in the episode; at test time, the learnt policy is applied to unseen data to generate the summary. However, learning such cross-input policies is very expensive because of the huge search space. Another issue is the *delayed rewards*: the ROUGE scores can be calculated only when the complete summary is generated. Such delayed rewards cause RL-based summarisers to take even longer time to converge. Although multiple techniques have been proposed to speed up the training of RL, e.g. memory replay [Mnih *et al.*, 2015], MIXER [Ranzato *et al.*, 2016] and imitation learning [Cheng *et al.*, 2018], training cross-input RL yet requires considerable time, data and parameter tuning.

On the other hand, input-specific RL (middle part in Fig. 1) neither requires parallel data (i.e. input documents and reference summaries for them) nor a reward oracle. Instead, it employs a handcrafted reward function with which the RL agent interacts in order to learn a policy specifically for the given input. By doing so, the size of the search space significantly decreases, which consequently diminishes the training time and computational resources. However, designing such a reward function is highly challenging as it should fit all inputs [Rioux *et al.*, 2014]. This explains the poor performance of input-specific RL summarisers [Ryang and Abekawa, 2012].

To tackle the problems faced by the above two RL-based summarisation paradigms, we propose a novel paradigm called *REward Learning for Input-Specific reinforcement learning (RELIS)*. Instead of learning a cross-input policy, RELIS learns a *cross-input reward oracle* at training time, and then uses the learnt reward to train an input-specific policy for each input at test time (bottom part in Fig. 1). RELIS is inspired by *inverse RL* [Abbeel and Ng, 2004], which requires a demonstrator to present optimal trajectories. Because such a demonstrator is hardly available in summarisation, RELIS leverages *Learning-to-Rank (L2R)* algorithms [Li, 2011] to approximate the ground-truth reward oracle from “weak supervisions”, e.g., numeric scores that indicate the quality of the summary or preferences over summary pairs.

\*Work performed while at UKP-TUDA.

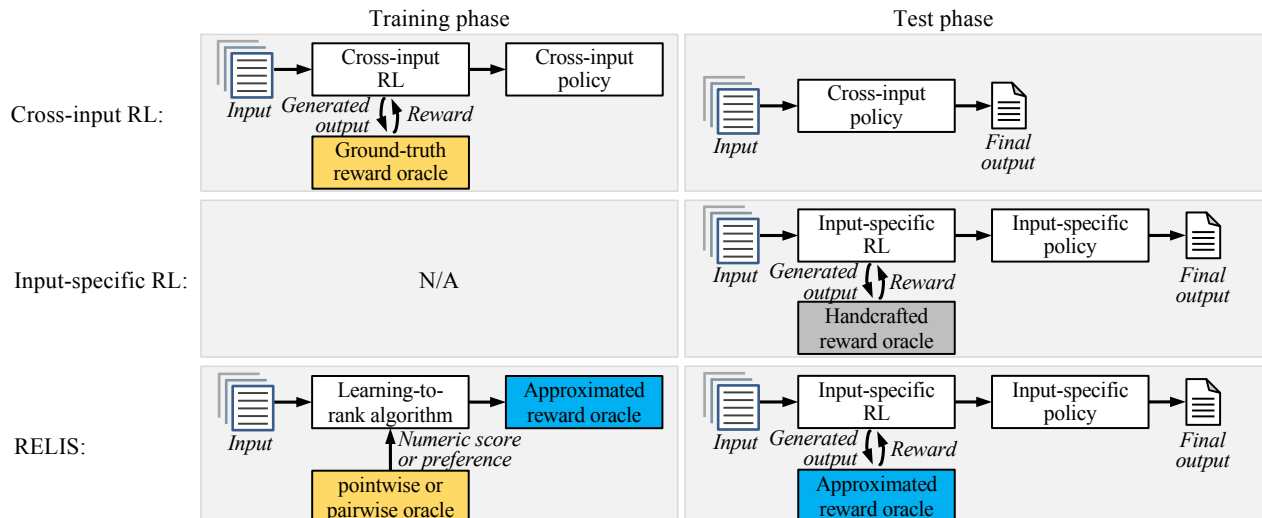


Figure 1: The workflows of cross-input RL (top), input-specific RL (middle) and RELIS (bottom). The ground-truth reward can be provided by humans or automatic metrics (e.g. BLEU or ROUGE) measuring the similarity between generated output text and the reference output.

Our contributions are threefold: **(i)** We propose RELIS (§2), a new RL-based summarisation framework that enjoys the strong performance of cross-input RL and the low computational cost of input-specific RL. **(ii)** Theoretically, we prove that by employing appropriate L2R and RL algorithms, RELIS guarantees to generate near-optimal outputs (§3). **(iii)** Empirically, we evaluate RELIS in multi-document extractive summarisation (§4). Compared to the state-of-the-art methods, RELIS provides comparable performance but requires much less the training time or data. Because the proof of RELIS is generic, we believe RELIS has the potential to be applied to other NLG tasks, e.g. translation and sentence simplification. Source code and supplementary material are available at <https://github.com/UKPLab/ijcai2019-relis>.

## 2 RELIS

We first formally define the summarisation task, and then detail the L2R and RL module of RELIS.

### 2.1 Extractive Summarisation

In line with Peyrard and Eckle-Kohler [2017], we formulate summarisation as a *discrete optimisation problem*. Let  $\mathcal{X}$  be the set of all possible input documents and  $\bar{\mathcal{X}}$  be the set of inputs available at training time. An input can be either a single document or a cluster of documents on the same topic. For input  $x \in \mathcal{X}$ , let  $\mathcal{Y}_x$  indicate the set of all extractive summaries for  $x$  that comply with the length requirement. Then, the task of summarisation is to map each input  $x$  to its best summary in  $\mathcal{Y}_x$  with respect to a *ranking function*  $\sigma_x: \mathcal{Y}_x \rightarrow \mathbb{N}$ . For a candidate summary  $y \in \mathcal{Y}_x$ ,  $\sigma_x$  returns the number of candidates in  $\mathcal{Y}_x$  that have equal or lower quality than  $y$ , including  $y$  itself. For example, if  $y$  is the highest-quality candidate in  $\mathcal{Y}_x$ , then  $\sigma_x(y) = |\mathcal{Y}_x|$ .  $\sigma_x$  can be obtained from human evaluators, automatic metrics (e.g. ROUGE) or heuristics measuring the quality of outputs. We denote the ground-truth ranking function on  $\mathcal{Y}_x$  by  $\sigma_x^*$ .

Given the above definition of summarisation, we define a summariser agent as a tuple  $(\sigma_x, M_x)$ , where  $M_x$  is an *optimisation model* for finding the (near-)optimal summary with respect to  $\sigma_x$ . In the cross-input paradigm, the RL agent learns a policy that solves the optimisation problem  $M_x$  for any  $x \in \mathcal{X}$  at training time. In RELIS, instead, an agent learns a ranking  $\hat{\sigma}_x^U: \mathcal{Y}_x \rightarrow \mathbb{N}$  at training time so that  $\hat{\sigma}_x^U$  is as “close” as possible to  $\sigma_x^*$  (“close” will be formally defined in §3). At test time, for each  $x' \in \mathcal{X} \setminus \bar{\mathcal{X}}$ , RELIS formulates  $M_{x'}$  as a *Markov Decision Process (MDP)* and learns an RL policy specifically for  $M_{x'}$ .

### 2.2 Learning to Rank (L2R)

L2R algorithms learn to reconstruct the ranking over objects from an *oracle* [Li, 2011]. L2R induces the approximated ranking  $\hat{\sigma}_x^U$  by learning a *utility function*  $U$  from the oracle, such that  $\hat{\sigma}_x^U(y_1) > \hat{\sigma}_x^U(y_2)$  if and only if  $U(y_1, x) > U(y_2, x)$ . L2R can learn from different types of oracles, including *pointwise* oracles that provide point-based scores for objects, *pairwise* oracles that provide preferences over pairs of objects, and *listwise* oracles that provide certain metrics (e.g. accuracy) for candidate rankings. Here, we focus on pointwise and pairwise oracles as humans reliably provide such judgements for short texts [Kreutzer *et al.*, 2018]. Function  $U$  can be learnt by any function approximation technique, e.g., neural networks.

In pointwise L2R, for every  $x \in \bar{\mathcal{X}}$ , we draw  $N$  sample summaries from  $\mathcal{Y}_x$  using some sampling strategy, e.g., random sampling without replacement. Then, we query their  $\sigma_x^*$  values from the oracle and use a regression algorithm to minimise the averaged mean squared error (MSE) between  $U$  and  $\sigma_x^*$ . Formally, the loss function is

$$\mathcal{L}^{\text{MSE}} = \frac{1}{N \cdot |\bar{\mathcal{X}}|} \sum_{x \in \bar{\mathcal{X}}} \sum_{i=1}^N (\sigma_x^*(y_i) - U(y_i, x))^2. \quad (1)$$

In pairwise L2R, for every  $x \in \overline{\mathcal{X}}$  we sample  $K$  pairs of summaries from  $\mathcal{Y}_x$  and then query the oracle about their preferences. We denote the collected preferences for input  $x$  by  $P_x = \{p_1(y_{1,1}, y_{1,2}), \dots, p_K(y_{K,1}, y_{K,2})\}$ , where  $y_{i,1}, y_{i,2} \in \mathcal{Y}_x$  are the candidate outputs presented to the oracle in the  $i^{\text{th}}$  iteration;  $p_i(y_{i,1}, y_{i,2})$  equals 1 if the oracle prefers  $y_{i,1}$  over  $y_{i,2}$ , and equals 0 otherwise. Different loss functions can be used to learn  $U$  from the preferences in  $P_x$ . First, we consider the cross-entropy loss:

$$\begin{aligned} \mathcal{L}^{\text{CE}} = & - \sum_{x \in \overline{\mathcal{X}}} \sum_{p_i(y_{i,1}, y_{i,2}) \in P_x} [p_i(y_{i,1}, y_{i,2}) \log \mathcal{P}(y_{i,1}, y_{i,2}) \\ & + p_i(y_{i,2}, y_{i,1}) \log \mathcal{P}(y_{i,2}, y_{i,1})], \end{aligned} \quad (2)$$

where  $\mathcal{P}(y_1, y_2) = (1 + \exp[U(y_2, x) - U(y_1, x)])^{-1}$ . An alternative is the margin ranking (a.k.a. pairwise hinge) loss:

$$\begin{aligned} \mathcal{L}^{\text{MR}} = & \frac{1}{N \cdot |\overline{\mathcal{X}}|} \sum_{x \in \overline{\mathcal{X}}} \sum_{p_i(y_{i,1}, y_{i,2}) \in P_x} \max[0, \\ & 1 - e_i \cdot (U(y_{i,1}, x) - U(y_{i,2}, x))], \end{aligned} \quad (3)$$

where  $e_i = 1$  if  $y_{i,1}$  is preferred over  $y_{i,2}$ , and  $e_i = -1$  otherwise. Additionally, we consider an improved margin ranking loss proposed by Agarwal and Collins [2010], which gives large penalties to mis-ranked pairs with wide margins:

$$\begin{aligned} \mathcal{L}^{\text{IM}} = & \frac{1}{N \cdot |\overline{\mathcal{X}}|} \sum_{x \in \overline{\mathcal{X}}} \sum_{p_i(y_{i,1}, y_{i,2}) \in P_x} \max[0, \\ & |\sigma_x^*(y_{i,1}) - \sigma_x^*(y_{i,2})| - e_i(U(y_{i,1}, x) - U(y_{i,2}, x))]. \end{aligned} \quad (4)$$

Note that the improved margin ranking loss is a mix of pointwise and pairwise L2R as it requires both  $\sigma_x^*$  and  $P_x$ .

### 2.3 Reinforcement Learning (RL)

RL amounts to algorithms that search for optimal solutions in *Markov Decision Processes (MDPs)*. We formulate the optimisation problems  $M_x$  for NLG tasks as episodic MDP [Ryang and Abekawa, 2012; Rioux *et al.*, 2014]. An episodic MDP is a tuple  $(S, A, P, R, T)$ , where  $S$  is a set of *states*,  $A$  is a set of *actions*,  $P: S \times A \rightarrow S$  is the *transition function* with  $P(s, a)$  giving the next state after performing action  $a$  in state  $s$ ,  $R: S \times A \rightarrow \mathbb{R}$  is the *reward function* with  $R(s, a)$  giving the immediate reward for performing action  $a$  in state  $s$  and  $T \subseteq S$  is a set of *terminal states* that mark the end of an episode. Furthermore, to formulate the delayed rewards in summarisation, we let  $R(s, a) = 0$  if  $P(s, a) \notin T$ , so as to ensure that non-zero rewards only appear in the last step of each episode.

In extractive summarisation, the components in  $M_x$  are defined as follows.  $S$  is the set of all possible states, i.e. permutations of sentences from  $x$ . Note that  $\mathcal{Y}_x$  is a subset of  $S$  because  $\mathcal{Y}_x$  only includes the summaries complying with the length requirement, but  $S$  includes all possible extractive summaries. Two types of action constitute  $A$ : *add* a new sentence from the input document cluster  $x$  to the current draft summary and *terminate* the generation.  $T$  includes all the over-length summaries and an *absorbing state*  $s_T$ . If action  $a$  is terminate,  $P(s, a) = s_T$  regardless of the current state  $s$ . The reward function  $R(s, a)$  returns  $\sigma_x(s)$  if action  $a$  is

terminate, a negative reward if the current state  $s$  is an over-length summary and 0 otherwise. We denote this MDP by  $M_x(\sigma_x)$  to highlight that it uses  $\sigma_x$  as rewards.

A *policy*  $\pi: S \times A \rightarrow [0, 1]$  defines how actions are selected:  $\pi(s, a)$  is the probability of selecting action  $a$  in state  $s$ . In the context of summarisation, we slightly abuse the notation by letting  $\pi(y|x)$  denote the probability of policy  $\pi$  generating  $y \in \mathcal{Y}_x$  given input  $x$ . We say a policy  $\pi$  is *proper* if  $\sum_{y \in \mathcal{Y}_x} \pi(y|x) = 1$ , i.e.  $\pi$  does not generate illegal outputs. Then, the expected reward of performing proper policy  $\pi$  is:

$$\mathcal{R}_{\sigma_x}(\pi|x) = \mathbb{E}_{y \sim \pi} \sigma_x(y) = \sum_{y \in \mathcal{Y}_x} \pi(y|x) \sigma_x(y). \quad (5)$$

The goal of RL is to find the optimal policy  $\pi^*$  that has the highest expected reward:  $\pi^* = \arg \max_{\pi} \mathcal{R}_{\sigma_x}(\pi|x)$ . Note that  $\pi_x^*$  is a probability distribution that assigns non-zero probabilities only to the optimal outputs for  $x$ , i.e.  $\pi^*(y|x) > 0$  only if  $\sigma_x(y) = |\mathcal{Y}_x|$ . Hence  $\mathcal{R}_{\sigma_x}(\pi^*|x) = |\mathcal{Y}_x|$ .

### 3 Proof of Convergence for RELIS

In this section, we prove that if the errors in both L2R and RL are bounded, the error of the final output of RELIS is a linear polynomial of the errors from the two steps. First, we define *convergent* L2R, which can learn an approximated ranking  $\hat{\sigma}_x^U$  “close” to the ground-truth ranking  $\sigma_x^*$ .

**Definition 1.** Let  $U$  be the utility function learnt by an L2R algorithm. For  $x \in \mathcal{X}$ , let  $\hat{\sigma}_x^U$  be the ranking derived from  $U$  on  $\mathcal{Y}_x$ . The L2R algorithm is  $\epsilon$ -convergent with respect to  $\sigma_x^*$  if there exists a non-negative integer  $\epsilon$  such that

$$\max_{y \in \mathcal{Y}_x} |\hat{\sigma}_x^U(y) - \sigma_x^*(y)| \leq \epsilon. \quad (6)$$

Some L2R algorithms are  $\epsilon$ -convergent under realistic conditions. For example, the quick-sort based L2R algorithm proposed by Maystre and Grossglauser [2017] is  $\epsilon$ -convergent with respect to pairwise oracles, when it obtains a sufficiently large number of preferences. Next, we define convergent RL algorithms that guarantee to learn near-optimal policies.

**Definition 2.** Let  $M_x(\sigma_x)$  denote the episodic MDP for an NLG task given input  $x \in \mathcal{X}$ , which uses a ranking function  $\sigma_x$  over  $\mathcal{Y}_x$  as its reward function. An RL algorithm is  $\delta$ -convergent for  $M_x(\sigma_x)$  if there exists  $\delta \in \mathbb{R}_{\geq 0}$  such that

$$\lim_{n \rightarrow \infty} \mathcal{R}_{\sigma_x}(\pi^*|x) - \mathcal{R}_{\sigma_x}(\pi_n|x) \leq \delta, \quad (7)$$

where  $\pi^*$  is the optimal policy for  $M_x(\sigma_x)$  and  $\pi_n$  is the policy of the RL algorithm after  $n$  episodes of learning.<sup>1</sup>

Many RL algorithms have been proven to be  $\delta$ -convergent, including value-based RL [Sutton *et al.*, 2009], actor-critic [Sutton *et al.*, 1999] and policy gradient [Fazel *et al.*, 2018].

Recall that for an input  $x$  from the test set  $\mathcal{X} \setminus \overline{\mathcal{X}}$ , our goal is to find the optimal policy  $\pi^*$  for the NLG tuple  $(\sigma_x^*, M_x(\sigma_x^*))$  (see §2). The theorem below shows that if RELIS uses convergent L2R and RL, its output is near-optimal, i.e. the error between the RELIS output and the optimal output is bounded.

<sup>1</sup>Since  $\pi^*$  is the optimal policy for  $\sigma_x$ ,  $\mathcal{R}_{\sigma_x}(\pi^*|x)$  is always higher than the expected reward of any other policy (cf. Eq. 5).

**Theorem 1.** For a test input  $x \in \mathcal{X} \setminus \overline{\mathcal{X}}$ , we denote the optimal NLG agent for  $x$  by  $(\sigma_x^*, M_x(\sigma_x^*))$ , and its RELIS approximation by  $(\hat{\sigma}_x^U, M_x(\hat{\sigma}_x^U))$ , where  $\hat{\sigma}_x^U$  is the approximation of  $\sigma_x^*$  learnt by an  $\epsilon$ -convergent L2R algorithm. Let  $\pi^*$  be the optimal policy for  $M_x(\sigma_x^*)$  and  $\hat{\pi}_n$  the proper policy for  $M_x(\hat{\sigma}_x^U)$  learnt by a  $\delta$ -convergent RL algorithm for  $n$  episodes. Then, as  $n$  approaches positive infinity,  $\mathcal{R}_{\sigma_x^*}(\hat{\pi}_n|x) \geq \mathcal{R}_{\sigma_x^*}(\pi^*|x) - \delta - \epsilon$ .

*Proof.* We denote the optimal policy for  $M_x(\hat{\sigma}_x^U)$  as  $\hat{\pi}^*$ . Hence from Eq. (7) we have:

$$\lim_{n \rightarrow \infty} \mathcal{R}_{\hat{\sigma}_x^U}(\hat{\pi}_n|x) \geq \mathcal{R}_{\hat{\sigma}_x^U}(\hat{\pi}^*|x) - \delta. \quad (8)$$

For any  $y \in \mathcal{Y}_x$ , from Eq. (6) we have

$$\hat{\sigma}_x^U(y) \leq \sigma_x^*(y) + \epsilon. \quad (9)$$

Using Eq. (9), we can bound  $\mathcal{R}_{\hat{\sigma}_x^U}(\hat{\pi}_n|x)$  with any positive  $n$  value (see Eq. (5)):

$$\begin{aligned} \mathcal{R}_{\hat{\sigma}_x^U}(\hat{\pi}_n|x) &= \sum_{x \in \mathcal{Y}_x} \hat{\pi}_n(y|x) \hat{\sigma}_x^U(y) \\ &\leq \sum_{x \in \mathcal{Y}_x} \hat{\pi}_n(y|x) [\sigma_x^*(y) + \epsilon] \\ &= \sum_{y \in \mathcal{Y}_x} \hat{\pi}_n(y|x) \sigma_x^*(y) + \epsilon \sum_{y \in \mathcal{Y}_x} \hat{\pi}_n(y|x) \\ &= \mathcal{R}_{\sigma_x^*}(\hat{\pi}_n|x) + \epsilon. \end{aligned} \quad (10)$$

Note that  $\sum_{y \in \mathcal{Y}_x} \hat{\pi}_n(y|x) = 1$  in Eq. (10) because  $\hat{\pi}_n$  is a proper policy (see §2). Combining Eq. (8) and (10), we get

$$\lim_{n \rightarrow \infty} \mathcal{R}_{\sigma_x^*}(\hat{\pi}_n|x) \geq \mathcal{R}_{\sigma_x^*}(\hat{\pi}^*|x) - \delta - \epsilon. \quad (11)$$

Since  $\hat{\pi}^*$  is the optimal policy for  $M_{\hat{\sigma}_x^U}$ , according to Eq. (5),  $\mathcal{R}_{\hat{\sigma}_x^U}(\hat{\pi}^*|x) = |\mathcal{Y}_x|$ . Similarly,  $\mathcal{R}_{\sigma_x^*}(\pi^*|x) = |\mathcal{Y}_x|$ . Hence we can replace  $\mathcal{R}_{\hat{\sigma}_x^U}(\hat{\pi}^*|x)$  in Eq. (11) with  $\mathcal{R}_{\sigma_x^*}(\pi^*|x)$  and obtain  $\lim_{n \rightarrow \infty} \mathcal{R}_{\sigma_x^*}(\hat{\pi}_n|x) \geq \mathcal{R}_{\sigma_x^*}(\pi^*|x) - \delta - \epsilon$ .  $\square$

## 4 Experimental Setup

**Task and datasets.** We evaluate RELIS for extractive multi-document summarisation on three benchmark datasets from the Document Understanding Conferences (DUC)<sup>2</sup> described in Table 1. Each dataset contains a set of document clusters. For each cluster, the target is to create a summary with at most 100 words. Each cluster is accompanied with several human-generated summaries, denoted as  $r_x$ , that we use for training and evaluation. To decide the best parameters, we perform 10-fold cross validation on DUC’01.

**Ground-truth reward oracle.** For  $x \in \overline{\mathcal{X}}$ , we use a ROUGE-based  $U^*$  to induce the ground-truth  $\sigma_x^*$ . We measure the quality of a summary  $y \in \mathcal{Y}_x$  by  $U^*(y, x) = R_1(y, r_x) + 2R_2(y, r_x)$ , where  $R_1$  and  $R_2$  denote the average ROUGE-1 and ROUGE-2 recall metrics, respectively.  $R_1$  and  $R_2$  are arguably the most widely used metrics to approximate human evaluation of summary quality. The  $R_2$  scores for optimal extractive summaries are usually half of their  $R_1$  scores [Gao *et al.*, 2018], so we multiply  $R_2$  by 2 to balance the impact of  $R_1$  and  $R_2$ . Next, we describe how we approximate the ground-truth reward.

<sup>2</sup><https://duc.nist.gov/>

Dataset	# Cluster	# Doc	# Sent/Cluster
DUC ’01	30	308	378
DUC ’02	59	567	271
DUC ’04	50	500	265

Table 1: Some statistics of the DUC datasets: the number of document clusters (# Cluster), the total number of documents (# Doc) and the average number of sentences per cluster (# Sent/Cluster).

### 4.1 Approximated Reward Oracle

Recall that the goal of L2R is to learn the utility function  $U$  using pointwise or pairwise preferences (see §2). In pointwise L2R, we sample  $N = 3000$  different summaries for each document cluster in  $x \in \overline{\mathcal{X}}$  by randomly selecting sentences from  $x$  until the length limit is reached. Using larger  $N$  does not significantly increase the quality of  $\hat{\sigma}_x^U$  but increases the training time. We use the ground-truth ranking  $\sigma_x^*$  over the  $|\overline{\mathcal{X}}| \cdot N$  samples to learn the utility function  $U$  by minimising the MSE loss from Eq. (1). In pairwise L2R, we randomly draw  $K = 10^5$  of the  $|\overline{\mathcal{X}}| \frac{N \cdot (N-1)}{2}$  possible pairs of the  $|\overline{\mathcal{X}}| \cdot N$  sampled summaries to compute  $U$  by minimising  $\mathcal{L}^{\text{CE}}$ ,  $\mathcal{L}^{\text{MR}}$  or  $\mathcal{L}^{\text{IM}}$  from Eq. (2) – (4). Preliminary results suggest that increasing  $K$  to  $10^6$  does not benefit the performance but significantly slows down the training, while decreasing it to  $10^4$  significantly harms the quality of  $\hat{\sigma}_x^U$ .

We use a linear model to approximate the utility function, i.e.  $U(y, x) = w \cdot \phi(y, x)$ , where  $\phi(y, x)$  encodes  $y$  for input  $x$  as a vector by concatenating the following features:

- The negative Jensen-Shannon divergence between unigrams, bigrams and named entities in  $y$  and  $x$ .
- The summary quality evaluation heuristics proposed by Rioux *et al.* [2014].
- TFIDF values averaged over all tokens and named entities in  $y$  [Peyrard and Gurevych, 2018].
- The average number of document clusters a word or named entity appears in.
- Rate of named entities from  $x$  appearing in  $y$ .
- The percentage of tokens in  $y$  that belong to some named entities in  $x$ .
- The redundancy feature proposed by Peyrard and Gurevych [2018], applied to unigrams, bigrams and named entities.

We additionally use two Convolutional Neural Network (CNN) architectures to generate auto-encoded features: StdCNN [Kim, 2014] and PriorSum [Cao *et al.*, 2015]. To train these auto-encoders, we follow Cao *et al.* [2015] and feed only the embeddings of the words of  $y$  into these models. We add a linear layer as the last layer to map the output vector of these models to a  $U$  value. Full settings of these models are in the supplementary material. For each  $x \in \mathcal{X} \setminus \overline{\mathcal{X}}$ , we measure the quality of the approximated ranking  $\hat{\sigma}_x^U$  by its ranking correlation to the ground-truth  $\sigma_x^*$ . We consider two ranking correlation metrics: Spearman’s  $\rho \in [-1, 1]$  and the Normalized Discounted Cumulative Gain (NDCG) on the top

	DUC'01		DUC'02		DUC'04	
	$\rho$	NDCG	$\rho$	NDCG	$\rho$	NDCG
ASRL	.176	.555	.131	.537	.145	.558
REAPER	.316	.638	.301	.639	.372	.701
JS	.549	.736	.525	.700	.570	.763
Our $\hat{\sigma}_x^U$	<b>.601</b>	<b>.764</b>	<b>.560</b>	<b>.727</b>	<b>.617</b>	<b>.802</b>

Table 2: The correlation of approximated and ground-truth ranking.  $\hat{\sigma}_x^U$  has significantly higher correlation over all other approaches.

1% items.  $\text{NDCG} \in [0, 1]$  puts more emphasis on the top elements with logarithmic decay weighting. For both metrics, higher values indicate stronger correlations.

We train the utility function  $U$  with the loss functions in Eq. (1)-(4) and we find that using different loss functions does not significantly<sup>3</sup> change performance (see supplementary material). However, the cross-entropy loss consistently results in marginally better performance than the others, and hence we use it throughout our experiments. We find that under all examined settings, the CNN-based features underperform the other features. Using the CNN-based features together with other features also worsen the quality of  $\hat{\sigma}_x^U$ . The reason is because both PriorSum and StdCNN only encode the summaries’ *document-independent features* [Cao *et al.*, 2015]. Encoding document-dependent features requires more sophisticated neural models [Wu and Hu, 2018; Narayan *et al.*, 2018] which require considerable time, data and parameter tuning. This undermines the benefits of RELIS. We leave the research for efficient reward representation learning for future work.

## 4.2 RELIS Setup

We test RELIS on all three DUC datasets as follows. In line with Cao *et al.* [2017] and Ren *et al.* [2018], we split train and test data in the “leave-one-out” manner so that documents from two datasets are used as the training set,  $\bar{\mathcal{X}}$ , and documents from the rest as the test set  $\mathcal{X} \setminus \bar{\mathcal{X}}$ . In each run in the “leave-one-out” experiments, we randomly select 30% data from the training set as the dev set, and select the model with the best performance on the dev set. We use Adam with initial learning rate  $10^{-2}$ . The number of epochs is 10 and batch size is 2. As for RL in RELIS, we use the same *temporal difference* algorithm as Rioux *et al.* [2014]. Full details of our parameter selection and results of using different loss functions are in the supplementary material.

## 5 Results

Table 2 compares the quality of our  $\hat{\sigma}_x^U$  with other widely used rewards for input-specific RL (see §4).  $\hat{\sigma}_x^U$  has significantly higher correlation to the ground-truth ranking compared with all other approaches, confirming that our proposed L2R method yields a superior reward oracle.

In Table 3, we compare RELIS with non-RL-based and RL-based summarisation systems. For non-RL-based systems, we report *ICSI* [Gillick and Favre, 2009] maximising

<sup>3</sup>In all experiments, we perform statistical significant test by two-tailed  $t$ -test and  $p$ -value  $< 0.01$ .

	DUC'01		DUC'02		DUC'04	
	$R_1$	$R_2$	$R_1$	$R_2$	$R_1$	$R_2$
ICSI	33.31	7.33	35.04	8.51	37.31	9.36
PriorSum	35.98	7.89	36.63	8.97	38.91	10.07
TCSum	<b>36.45</b>	7.66	36.90	8.61	38.27	9.66
TCSum <sup>-</sup>	33.45	6.07	34.02	7.39	35.66	8.66
SRSum	36.04	8.44	<b>38.93</b>	<b>10.29</b>	39.29	10.70
DeepTD	28.74	5.95	31.63	7.09	33.57	7.96
REAPER	32.43	6.84	35.03	8.11	37.22	8.64
RELIS	34.73	<b>8.66</b>	37.11	9.12	<b>39.34</b>	<b>10.73</b>

Table 3: Results of non-RL (top), cross-input (DeepTD) and input-specific (REAPER) RL approaches (middle) compared with RELIS.

the bigram overlap of summary and input using integer linear programming, *PriorSum* [Cao *et al.*, 2015] learning sentence quality with CNNs, *TCSum* [Cao *et al.*, 2017] employing text classification of the input documents, the variant of TCSum without the text classification pre-training (TCSum<sup>-</sup>) and *SRSum* [Ren *et al.*, 2018], which learns sentence relations with both word- and sentence-level attentive neural networks to estimate salience.

For RL-based systems, we re-implement REAPER [Rioux *et al.*, 2014] as an input-specific RL, and DeepTD as a cross-input RL. DeepTD is adapted from the DQN-based RL summariser [Yao *et al.*, 2018] and is trained by taking  $U^*$  as the rewards (see §4). It uses InferSent to represent summaries. To improve the efficiency and performance of DeepTD, we use *memory replay* and *periodic update*. Further details and the algorithm of DeepTD are in the supplementary material.

RELIS significantly outperforms the other RL-based systems. Note that RELIS and REAPER use the identical RL algorithm for input-specific policy learning; hence the improvement of RELIS is due to the higher quality of the L2R-learned reward  $\hat{\sigma}_x^U$ . RELIS outperforms DeepTD because training cross-input policies requires much more data than available in the DUC datasets. At the same time, RELIS performs on par with neural-based TCSum and SRSum, while it requires significantly less data and time to train, as shown next.

We run RELIS, SRSum, DeepTD and REAPER on the same workstation with a 4-core CPU, 8 GB memory and no GPUs. Table 4 compares their average training and test time for each document cluster. RELIS reduces the training time of SRSum by two orders of magnitude. At test time, RELIS takes reasonable time to train the input-specific policy, and we believe that it can be further reduced by using more efficient RL algorithms and employing techniques like memory replay or reward shaping.

Unlike TCSum, RELIS requires no additional training data: TCSum uses 1.8 million news articles from New York Times and their category annotations (e.g. Health, Politics, Business) for training. It is worth noting that without using such a massive extra data for the text classification step, the performance of TCSum significantly drops (see TCSum<sup>-</sup> in Table 3). The training time of TCSum is unlikely to be shorter than RELIS since TCSum requires to train a CNN-based text classifier before training a CNN-based sentence selector.

	SRSUM	DeepTD	REAPER	RELIS
Training	810 s	1,560 s	N/A	2 s
Test	7 s	4 s	31 s	34 s

Table 4: Averaged training and test time for each document cluster. Note that REAPER does not have a training phase.

To summarise, RELIS significantly outperforms RL-based models, and it yields competitive performance compared with the state-of-the-art neural summarisers, with the benefit of needing less training time and data.

## 6 Discussion & Related Work

RELIS is proposed as a RL-based summarisation paradigm, but it can be applied to other NLG tasks where RL is widely used, e.g., translation, sentence simplification and dialogue generation. For example, to apply RELIS to translation, we just let  $\mathcal{X}$  (see §2) be the set of texts in the source language and let  $\mathcal{Y}_x$  be the set of all possible translations in the target language for input  $x$ ; and the error bound of RELIS (Theorem 1) still holds as it is indifferent to the contents in  $\mathcal{X}$  and  $\mathcal{Y}_x$ . Hence, in this section, we discuss the related work in the context of NLG in general.

*Reward learning* recently receives increasing interests from the machine learning community [Ibarz *et al.*, 2018; Zheng *et al.*, 2018], but it is largely overlooked in NLG until now. Unlike classic RL applications, e.g. robotics and games, where rewards are provided or easy to design, NLG tasks lack strong metrics to measure the quality of the output. Well-known rewards, e.g. ROUGE, are criticised for their low correlation with human evaluations [Chaganty *et al.*, 2018]. Recent work suggests that improving the reward function boosts the performance of RL-based NLG [Kryscinski *et al.*, 2018].

Besides using metrics such as BLEU and ROUGE as rewards to train RL-based NLG, novel rewards have been designed. Kryscinski *et al.* [2018] propose a new reward function for encouraging RL-based summarisers to prefer *novel* words in abstractive summaries. For sentence simplification, Zhang and Lapata [2017] propose a reward function measuring the simplicity, relevance and grammatical correctness of candidate outputs. For machine translation, Nguyen *et al.* [2017] propose a simulated user, which simulates human’s ratings on translations, and they use the simulated user to provide rewards for an RL-based translator. However, all rewards discussed above require reference outputs, unlike our learnt reward function which can be used to provide rewards at test time when no reference outputs are available.

Wu and Hu [2018] and Bosselut *et al.* [2018] recently propose to use a large volume of unlabelled data to learn a scorer measuring the *discourse coherence degree* of sentences, and use the scorer as rewards to train cross-input RL. We go beyond their work by proving the soundness of the combination of reward learning and RL, and training input-specific instead of cross-input policies so as to reduce training time and data.

RL-based *interactive NLG* methods elicit human feedback as rewards. Kreutzer *et al.* [2018] elicit pointwise and pairwise feedback on candidate translations to train a cross-input

policy. Gao *et al.* [2018] propose using active learning to select appropriate candidate summary pairs and acquire human preferences to improve a heuristics-based reward. However, these methods require much feedback to bring satisfactory results, e.g., at least 50 summary pairs to yield significant improvements [Gao *et al.*, 2018]. RELIS can be used as a pre-training stage for interactive methods, as RELIS first learns a high-quality cross-input reward function and then uses the interactive NLG techniques to elicit a small amount of feedback to fine-tune a user- and input-specific reward function, so as to generate higher-quality and personalised results.

## 7 Conclusion

We propose a novel RL paradigm called RELIS, which learns a reward function from a reward oracle with learning-to-rank (L2R) algorithms at training time, and then uses the learnt reward to train input-specific RL policies at test time. Compared with the widely employed cross-input RL-based summarisation approaches, RELIS avoids the expensive learning of cross-input policies but, instead, efficiently performs L2R and input-specific RL learning. Moreover, RELIS avoids the arduous reward design required in input-specific RL-based summarisation approaches. We prove that, with proper L2R and RL algorithms, RELIS guarantees to produce near-optimal outputs. Our experiments show that even with linear L2R and standard RL algorithms, RELIS yields performance on par with the state of the art while only requiring a small fraction of data and time to train. Our work lays the theoretical foundation for reward learning in NLG, and we hope it will encourage further research in this direction.

## Acknowledgements

This work has been supported by the German Research Foundation (DFG), as part of the QA-EduInf project (GU 798/18-1 and RI 803/12-1) and through the German-Israeli Project Cooperation (DIP, DA 1600/1-1 and GU 798/17-1).

## References

- [Abbeel and Ng, 2004] Pieter Abbeel and Andrew Y. Ng. Apprenticeship learning via inverse reinforcement learning. In *ICML*, 2004.
- [Agarwal and Collins, 2010] Shivani Agarwal and Michael Collins. Maximum margin ranking algorithms for information retrieval. In *ECIR*, pages 332–343, 2010.
- [Bosselut *et al.*, 2018] Antoine Bosselut, Asli Çelikyilmaz, Xiaodong He, Jianfeng Gao, Po-Sen Huang, and Yejin Choi. Discourse-aware neural rewards for coherent text generation. In *NAACL-HLT*, pages 173–184, 2018.
- [Cao *et al.*, 2015] Ziqiang Cao, Furu Wei, Sujian Li, Wenjie Li, Ming Zhou, and Houfeng Wang. Learning summary prior representation for extractive summarization. In *ACL*, pages 829–833, 2015.
- [Cao *et al.*, 2017] Ziqiang Cao, Wenjie Li, Sujian Li, and Furu Wei. Improving multi-document summarization via text classification. In *AAAI*, pages 3053–3059, 2017.

- [Chaganty *et al.*, 2018] Arun Tejasvi Chaganty, Stephen Mussmann, and Percy Liang. The price of debiasing automatic metrics in natural language evaluation. In *ACL*, pages 643–653, 2018.
- [Cheng *et al.*, 2018] Ching-An Cheng, Xinyan Yan, Nolan Wagener, and Byron Boots. Fast policy learning through imitation and reinforcement. In *UAI*, pages 845–855, 2018.
- [Fazel *et al.*, 2018] Maryam Fazel, Rong Ge, Sham Kakade, and Mehran Mesbahi. Global convergence of policy gradient methods for the linear quadratic regulator. In *ICML*, pages 1466–1475, 2018.
- [Gao *et al.*, 2018] Yang Gao, Christian M. Meyer, and Iryna Gurevych. APRIL: Interactively learning to summarise by combining active preference learning and reinforcement learning. In *EMNLP*, pages 4120–4130, 2018.
- [Gillick and Favre, 2009] Dan Gillick and Benoit Favre. A scalable global model for summarization. In *ILP*, pages 10–18. Association for Computational Linguistics, 2009.
- [Ibarz *et al.*, 2018] Borja Ibarz, Jan Leike, Tobias Pohlen, Geoffrey Irving, Shane Legg, and Dario Amodei. Reward learning from human preferences and demonstrations in atari. In *NeurIPS*, pages 8022–8034, 2018.
- [Kim, 2014] Yoon Kim. Convolutional neural networks for sentence classification. In *EMNLP*, pages 1746–1751, 2014.
- [Kreutzer *et al.*, 2018] Julia Kreutzer, Joshua Uyheng, and Stefan Riezler. Reliability and learnability of human bandit feedback for sequence-to-sequence reinforcement learning. In *ACL*, pages 1777–1788, 2018.
- [Kryscinski *et al.*, 2018] Wojciech Kryscinski, Romain Paulus, Caiming Xiong, and Richard Socher. Improving abstraction in text summarization. In *EMNLP*, pages 1808–1817, 2018.
- [Li, 2011] Hang Li. A short introduction to learning to rank. *IEICE Transactions*, 94-D(10):1854–1862, 2011.
- [Lin, 2004] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *ACL Workshop “Text Summarization Branches Out”*, pages 74–81, 2004.
- [Maystre and Grossglauser, 2017] Lucas Maystre and Matthias Grossglauser. Just sort it! A simple and effective approach to active preference learning. In *ICML*, pages 2344–2353, 2017.
- [Mnih *et al.*, 2015] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellefleur, Alex Graves, Martin A. Riedmiller, Andreas Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [Narayan *et al.*, 2018] Shisha Narayan, Shay B. Cohen, and Mirella Lapata. Ranking sentences for extractive summarization with reinforcement learning. In *NAACL-HLT*, pages 1747–1759, 2018.
- [Nguyen *et al.*, 2017] Khanh Nguyen, Hal Daumé III, and Jordan L. Boyd-Graber. Reinforcement learning for bandit neural machine translation with simulated human feedback. In *EMNLP*, pages 1465–1475, 2017.
- [Peyrard and Eckle-Kohler, 2017] Maxime Peyrard and Judith Eckle-Kohler. A principled framework for evaluating summarizers: Comparing models of summary quality against human judgments. In *ACL*, pages 26–31, 2017.
- [Peyrard and Gurevych, 2018] Maxime Peyrard and Iryna Gurevych. Objective function learning to match human judgements for optimization-based summarization. In *NAACL-HLT*, pages 654–660, 2018.
- [Ranzato *et al.*, 2016] Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. Sequence level training with recurrent neural networks. In *ICLR*, 2016.
- [Ren *et al.*, 2018] Pengjie Ren, Zhumin Chen, Zhaochun Ren, Furu Wei, Liqiang Nie, Jun Ma, and Maarten de Rijke. Sentence relations for extractive summarization with deep neural networks. *ACM Trans. Inf. Syst.*, 36(4):39:1–39:32, 2018.
- [Rioux *et al.*, 2014] Cody Rioux, Sadid A. Hasan, and Yilias Chali. Fear the REAPER: A system for automatic multi-document summarization with reinforcement learning. In *EMNLP*, pages 681–690, 2014.
- [Ryang and Abekawa, 2012] Seonggi Ryang and Takeshi Abekawa. Framework of automatic text summarization using reinforcement learning. In *EMNLP/CoNLL*, pages 256–265, 2012.
- [Sutton *et al.*, 1999] Richard S. Sutton, David A. McAllester, Satinder P. Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *NIPS*, pages 1057–1063, 1999.
- [Sutton *et al.*, 2009] Richard S. Sutton, Hamid Reza Maei, Doina Precup, Shalabh Bhatnagar, David Silver, Csaba Szepesvári, and Eric Wiewiora. Fast gradient-descent methods for temporal-difference learning with linear function approximation. In *ICML*, pages 993–1000, 2009.
- [Williams, 1992] Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8:229–256, 1992.
- [Wu and Hu, 2018] Yuxiang Wu and Baotian Hu. Learning to extract coherent summary via deep reinforcement learning. In *AAAI*, pages 5602–5609, 2018.
- [Yao *et al.*, 2018] Kaichun Yao, Libo Zhang, Tiejian Luo, and Yanjun Wu. Deep reinforcement learning for extractive document summarization. *Neurocomputing*, 284:52–62, 2018.
- [Zhang and Lapata, 2017] Xingxing Zhang and Mirella Lapata. Sentence simplification with deep reinforcement learning. In *EMNLP*, pages 584–594, 2017.
- [Zheng *et al.*, 2018] Zeyu Zheng, Junhyuk Oh, and Satinder Singh. On learning intrinsic rewards for policy gradient methods. In *NeurIPS*, pages 4649–4659, 2018.