RST-LoRA: A Discourse-Aware Low-Rank Adaptation for Long Document Abstractive Summarization

Department of Computer Science Department of Language Science and Technology Saarland Informatics Campus, Saarland University, Germany dongqi.me@gmail.com



UNIVERSITÄT DES SAARLANDES



European Research Council Established by the European Commissio

erc



Dongqi Pu, Vera Demberg

RST-LoRA improves long document summarization by integrating methods.

TL;DR

rhetorical structure theory into the LoRA model, outperforming previous

Motivation

- Why do we need **low-rank** approximation?
- Why do we need **discourse** knowledge?

- Why do we need low-rank approximation?



Vanilla full-parameter fine-tuning (FFT)

 $N \times M$





- Why do we need low-rank approximation?



Motivation

- Why do we need low-rank approximation?

 - Only 0.01–1% of the parameters, PEFTs \approx FFT



FFT vs PEFT

 \checkmark



Motivation

- Why do we need discourse knowledge?
 - Challenges in PEFTs
 - Latent text relations
 - Importance level of different sentences

EPFTs are not driven or guided by discourse knowledge during the training phase, as this is not explicitly present in the input data.

Ghazvininejad et al. (2022); Zhao et al. (2023)

Reason

RST Prerequisite

- Rhetorical Structure Theory (RST) is helpful for determining:
 - Which sentences should or should not be included in the summary
 - Sentences relations
 - Discourse importance level

RST Prerequisite Cause-Effect EDU1 Elaboration

- **EDU1** is the most pivotal component
- **EDU2** provides information for **EDU3**
 - It is not a problem to delete EDU2
 - It is still fine to delete both EDU2 and 3

Figure 1: An example of RST tree: [Utilizing discourse structure to enhance text summarization is beneficial.]^{EDU1} [This technique can be used to identify key ideas and capture often overlooked nuances.]^{EDU2} [Accurate capture of these complex structures facilitates the generation of good summaries.]^{EDU3}

EDU2

 \rightarrow Nucleus

Satellite —

EDU3

- RST Distribution
- RST-Aware Injection

Our Method

- RST Distribution
 - Each point indicates the probability value $p(edu_i, edu_i, r_k) \in [0,1] \subseteq \mathbb{R}$ that edu_i is the nucleus of edu_i with discourse relation r_{k} . (Pu et al., ACL 2023)
 - We average and merge the y-axis of the matrix, and the merged value $c(edu_i, edu_i, r_k)$ is called the importance index of edu_i with relation r_k .

Our Method





- **RST Distribution (4 variants)**
 - RST_{wo}^b : Binary, label-agnostic representation (1 or 0)
 - RST_w^b : Binary distribution with relation labels
 - RST^p_{wo} : Label-omitted probabilistic representation
 - RST_w^p : Most fine-grained representation with relation types and probabilities

Our Method



- Experimental Settings
 - Datasets
 - Parser
 - Metrics
 - Training and Inference

• Datasets

- Multi-LexSum (ML, Shen et al., 2022)
- eLife (Goldsack et al., 2022)
- BookSum Chapter (BC, Kryscinski et al., 2022)

• Parser

- DMRST (Liu et al., 2020, 2021).

Extracting probabilities and type labels from final logits layer

• Metrics

- Rouge-Lsum (RLsum) (Lin, 2004)
- BERTScore (Zhang et al., 2020)
- METEOR (Banerjee and Lavie, 2005)
- sacreBLEU (Post, 2018)
- NIST (Lin and Hovy, 2003)

F1 scores of Rouge-1 (R1), Rouge-2 (R2), Rouge-L (RL), and

• Training and Inference

- Backbones

 - Vicuna13B-16k (Zheng et al., 2023) *GPT*
- Baselines
 - Backbones w/ FFT
 - Backbones w/ LoRA
 - GPT-4 (in-context learning)
- Other SOTAs



Main Results

LoRA vs. FFT: Comparable, more efficient

RST variant performance



 RST_w^p -LoRA: Best performance

Our best model vs GPT-4 and SOTAs



Hallucination Checking

SummaC testing: 0-1 score range

- GPT-4: Weakest consistency
- **RST enhances LoRA:** Reduces hallucinations

GPT4 - ML FFT - ML LoRA - ML RST $^{p}_{W}$ - ML GPT4 - eLife FFT - eLife LoRA - eLife RST $^{p}_{W}$ - eLife GPT4 - BC FFT - BC LoRA - BC

Candidate





Impact of Parser Capability



Impact of Random Masking on the Parser

Impact of Different Rank r



r = 8 is a trade-off point between performance gain and computational cost

- Human evaluation: BookSum, 10 instances
- Evaluators: CL/CS Graduate candidates, blind test
- RST_w^p -LoRA: Highest neural model performance



Human Evaluation



- **GPT-4** self-evaluation: Lowest scores to own answers
- RST_w^p -LoRA: more closer to the quality of humangenerated summaries





Conclusion

- model.
- Discourse uncertainty and relation labels are complementarily.
- Our model outperforms current SOTA models in specific evaluation metrics.

A method for injecting discourse knowledge into the training of LoRA

More Info

- Code: <u>https://dongqi.me/projects/RST-LoRA</u>
- Questions: dongqi.me@gmail.com



Thanks for listening Q&A



This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 Research and Innovation Programme (Grant Agreement No. 948878)



Established by the European Commission