# Improving Faithfulness by Augmenting Negative Summaries from Fake Documents

**Tianshu Wang[1], Faisal Ladhak[2], Esin Durmus[3], He He[1]**
[1]New York University, [2] Columbia University, [3] Stanford University
{tw2112, hhe}@nyu.edu, faisal@cs.columbia.edu, esdurmus@stanford.edu

## Abstract

Current abstractive summarization systems tend to hallucinate content that is unfaithful to the source document, posing a risk of misinformation. To mitigate hallucination, we must teach the model to distinguish hallucinated summaries from faithful ones. However, the commonly used maximum likelihood training does not disentangle factual errors from other model errors. To address this issue, we propose a back-translation-style approach to augment negative samples that mimic factual errors made by the model. Specifically, we train an *elaboration* model that generates hallucinated documents given the reference summaries, and then generates negative summaries from the fake documents. We incorporate the negative samples into training through a controlled generator, which produces faithful/unfaithful summaries conditioned on the control codes. Additionally, we find that adding textual entailment data through multitasking further boosts the performance. Experiments on three datasets (XSum, GigaWord, and WikiHow) show that our method consistently improves faithfulness without sacrificing informativeness according to both human and automatic evaluation.[1]

## 1 Introduction

Despite the fast progress in fluency and coherence of text summarization systems, a common challenge is that the generated summaries are often unfaithful to the source document, containing hallucinated, non-factual content (Cao et al., 2018; Falke et al., 2019, *inter alia*). Current summarization models are usually trained by maximum likelihood estimation (MLE), where unfaithful and faithful summaries are penalized equally if they both deviate from the reference. As a result, when
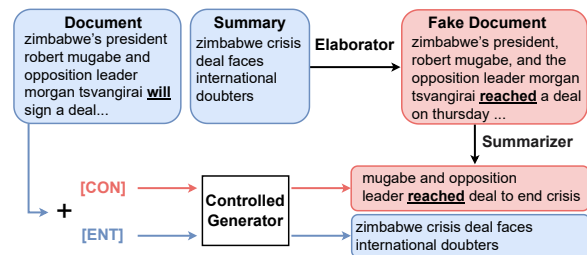


Figure 1: Overview of CoFE. The original and fabricated document-summary pairs are shown in blue and red respectively. The trained elaborator first generates fake documents from the summary. Then, the summarizer generates summaries from the fake documents, which are likely to contain hallucinated information (underlined). A controlled generator is then trained to produce the original (faithful) and the fabricated (unfaithful) summaries depending on the control codes.

the model fails to imitate the reference, it is likely to "over-generalize" and produce hallucinated content.

In this work, we address the issue by explicitly teaching the model to discriminate between positive (groundtruth) and negative (unfaithful) summaries. The key challenge is to generate realistic negative samples. Existing work on negative data augmentation mainly focuses on corrupting the reference (e.g., replacing entities) or sampling low-probability model outputs (Cao and Wang, 2021; Kryscinski et al., 2020; Kang and Hashimoto, 2020). However, the synthetic data often does not resemble actual hallucinations from the model (Goyal and Durrett, 2021) and many methods rely on external tools such as NER taggers.

To generate unfaithful summaries, we propose a simple method inspired by back-translation (Sennrich et al., 2016) (Fig. 1). Specifically, we first generate fake documents using an *elaboration* model that is trained to produce a document given the summary. We then generate summaries from

---

11913

the fake documents, which are assumed to be unfaithful since they are likely to contain hallucinated information in the fake documents. Given the reference summaries and the augmented negative samples, we train a controlled generation model that generates either faithful or unfaithful summaries conditioned on a faithfulness control code. At inference time, we control the model to generate only faithful summaries. We call our approach CoFE (**Co**ntrolled **F**aithfulness via **E**laboration). The controlled generation framework allows us to incorporate additional data easily: jointly training on natural language inference (NLI) datasets to generate entailed (faithful) and non-entailed (unfaithful) hypotheses further improves the result.

We evaluate CoFE on three summarization datasets: XSum (Narayan et al., 2018), GigaWord (Graff et al., 2003), and WikiHow (Koupaee and Wang, 2018). Both automatic metrics and human evaluation show that our method consistently outperforms previous methods in terms of faithfulness and content similarity to the reference, without sacrificing abstractiveness (Ladhak et al., 2022).

## 2 Approach

To learn a summarization model, the commonly used MLE aims to imitate the reference and does not distinguish different types of errors, thus the model may be misaligned with the desired behavior in downstream applications. For example, a faithful summary missing a detail would be preferred over a summary with hallucinated details, even if both have low likelihood under the data distribution. Therefore, additional inductive bias is needed to specify what unfaithful summaries are. Therefore, we augment negative examples and jointly model the distributions of both faithful and unfaithful summaries. At decoding time, we generate the most likely *faithful* summary.

**Negative data augmentation.** The key challenge in generating negative summaries is to simulate actual model errors. Prior approaches largely focus on named entities errors. However, different domains exhibit diverse hallucination errors (Goyal and Durrett, 2021); in addition, certain domains may not contain entities that can be easily detected by off-the-shelf taggers (e.g., stories or instructions). Our key insight is that the reverse summarization process—expanding a summary into a document—requires the model to hallucinate details, thus provides a domain-general way to pro-

duce unfaithful information. Instead of manipulating the reference summary directly, we expand it into a fake document, and generate negative summaries from it using the summarization model.

More formally, given a set of document-summary pairs $(x, y)$, we train a backward elaboration model $p_{\text{back}}(x \mid y)$ as well as a forward summarization model $p_{\text{for}}(y \mid x)$. Then, given a reference summary $y$, we first generate a fake document $\hat{x}$ from $p_{\text{back}}$, then generate the negative sample $y_{\text{neg}}$ from $\hat{x}$ using $p_{\text{for}}$, forming a pair of positive and negative samples $(x, y)$ and $(x, y_{\text{neg}})$. To avoid data leakage (i.e. training models and generating summaries on the same data), we split the training data into $K$ folds; the negative examples in each fold are generated by elaboration and summarization models trained on the rest $K - 1$ folds. We use $K = 5$ in the experiments.

**Controlled generation.** Given the positive and negative samples, we would like the model to learn to discriminate faithful summaries from unfaithful ones. Inspired by controlled generation methods (Keskar et al., 2019), we train the model to generate faithful or unfaithful summaries conditioned on a control code. In practice, we prepend a prefix at the beginning of the document (`[ENT]` for positive examples and `[CON]` for negative examples). At inference time, we always prepend `[ENT]` to generate faithful summaries.

**Training.** Our training data consists of positive examples (i.e. the original dataset) and generated negative samples, marked with different prefixes. Let $\mathcal{L}_{\text{pos}}, \mathcal{L}_{\text{neg}}$ denote negative log-likelihood (NLL) losses on the positive and negative examples. We use a multitasking loss that is a weighted sum of the two losses to balance the contribution from different types of examples: $\mathcal{L} = \mathcal{L}_{\text{pos}} + \lambda_1 \mathcal{L}_{\text{neg}}$ .

**Adding NLI datasets.** We hypothesize that incorporating NLI data through multitasking would transfer knowledge of entailment to the generator, helping it better model faithful and unfaithful summaries. The NLI sentence pairs can be naturally incorporated into controlled generation. Specifically, given the premise as input, we generate entailed and non-entailed hypotheses with control codes `[ENT]` and `[CON]`, respectively. With the additional NLI data, The loss function becomes: $\mathcal{L} = \mathcal{L}_{\text{pos}} + \lambda_1 \mathcal{L}_{\text{neg}} + \lambda_2 \mathcal{L}_{\text{NLI}}$ , where $\mathcal{L}_{\text{NLI}}$ denotes the NLL loss on the auxiliary NLI examples.

## 3 Experiments

**Datasets.** We evaluate our approach on 3 datasets,including: (i) **XSum** (Narayan-Chen et al., 2019), a dataset of BBC news articles paired with one-sentence summaries; (ii) **GigaWord** (Rush et al., 2015), a headline generation dataset compiled from the GigaWord corpus (Graff et al., 2003); and (iii) **WikiHow** (Koupaee and Wang, 2018), a dataset of how-to articles compiled from WikiHow.com, each paired with paragraph headlines as the summary. For the auxiliary NLI data, we use **SNLI** (Bowman et al., 2015) and **MultiNLI** (Williams et al., 2018), both containing pairs of premise and hypothesis sentences.

**Baselines.** We compare with three baselines: (i) maximum likelihood estimation (**MLE**); (ii) Loss Truncation (**LT**) (Kang and Hashimoto, 2020) that adaptively removes high-loss examples, which are assumed to be noisy/unfaithful; and (iii) **CLIFF** (Cao and Wang, 2021), a contrastive learning method based on generated negative samples.[2]

**Implementation.** All generation models (including the baselines) are fine-tuned BART-large models (Lewis et al., 2019). We train all CoFE models using Fairseq (Ott et al., 2019) with a learning rate of 3e-5. For decoding, we use beam search with a beam size of 6. We train the elaborators using the same model and learning hyperparameters. We generate one negative sample per document using beam search except for WikiHow where we use top-5 sampling.[3] To ensure that the negative summaries are different from the references, we further remove the top 10% summaries ranked by their edit distances to the reference. To train the controlled generator, we set coefficients ($\lambda_1, \lambda_2$) of the loss terms such that the reweighted number of examples in the original dataset, the negative samples, and optionally the NLI datasets have the ratio $1 : 0.5 : 0.5$. Details for other baselines are given in Appendix B.

**Metrics.** A good summary must cover important content, be faithful to the document, and be succinct. We evaluate the generated summaries from the following aspects. (1) *Content selection*. We use similarity to the reference as a proxy measure, and report ROUGE (Lin, 2004) and BertScore (Zhang et al., 2020). (2) *Faithfulness*. For automatic evaluation, we use QuestEval (Scialom et al., 2021), a QA-based metric, which shows better correlation with human judgment on system ranking in our preliminary experiments. We perform human evaluation on 100 randomly selected examples from each dataset. Given a document with the generated summaries from all systems (including the references), we ask annotators from Amazon Mechanical Turk to evaluate whether each summary is supported by the document. Each output is evaluated by 3 annotators. If two or more annotators vote "supported", then we consider the output faithful. The evaluation interface is described in Appendix A. (3) *Extractiveness*. **?** show that it is important to measure the extractiveness of the summaries to determine whether a method improves faithfulness mainly by copying from the document. Therefore, we also report *coverage* and *density* that measure the percentage of the words and the average length of text spans copied from the document (Grusky et al., 2018).

**Results.** Table 1 shows our main results. CoFE outperforms the baselines in human evaluated faithfulness accuracy on 2 out of the 3 datasets. On GigaWord, LT performs the best but it also incurs the largest drop in ROUGE and BertScore and more copying. CLIFF is good at fixing entity errors, but it has less advantage on datasets like WikiHow that contain fewer entities detectable by off-the-shelf taggers. On average, CoFE is less extractive than CLIFF and LT, indicating that our faithfulness improvements are not simply due to more copying. Finally, we find that adding NLI brings a marginal improvement on top of our negative samples.

**Are generated negative summaries really unfaithful?** Our method relies on the assumption that the elaboration of summaries introduces hallucinations, which results in unfaithful summaries. To verify this, we assess whether our generated negative samples are true negatives. Specifically, we evaluate the faithfulness of the negative summaries generated by our method and CLIFF on 100 randomly sampled documents from each dataset. In Table 2, we report the QuestEval scores and human-annotated faithfulness scores (following the same procedure described in Metrics). As a sanity check, the faithfulness scores of the negative samples are

---

[2]For CLIFF, we use SysLowCon which is reported to be the best amongst their methods for negative sample generation.

[3]WikiHow has very short summaries and we found it easy to generate the original references, thus we use sampling to increase diversity.

| Dataset | Method | Ref. Similarity (↑) | | Faithfulness (↑) | | Extractiveness (↓) | |
|---|---|---|---|---|---|---|---|
| | | RL | BS | Human Acc | QuestEval | Coverage | Density |
| XSum | MLE | **37.21** | 45.36 | 64% / 192 | 45.22 | 0.7596 | 1.6986 |
| | LT | 35.77 | 47.39 | 61% / 188 | 45.26 | 0.7564 | 1.7473 |
| | CLIFF | 36.41 | 52.78 | 68% / 192 | 45.48 | 0.7670 | 1.6904 |
| | CoFE | 36.38 | 52.09 | 68% / 194 | 45.54 | 0.7534 | 1.6460 |
| | CoFE +NLI | 36.98 | **52.90** | **70% / 196** | **45.98** | **0.7528** | **1.5961** |
| | CLIFF(CoFE data) | 36.06 | 52.35 | - | 45.33 | 0.7634 | 1.6703 |
| | CoFE(CLIFF data) | 36.73 | 52.42 | - | 45.23 | 0.7551 | 1.6207 |
| GigaWord | MLE | 33.95 | 27.77 | 70% / 206 | 43.80 | **0.7302** | **1.9415** |
| | LT | 34.22 | 26.35 | 76% / 204 | **45.58** | 0.8026 | 2.7106 |
| | CLIFF | **35.59** | **30.78** | 73% / 201 | 43.98 | 0.7406 | 2.1100 |
| | CoFE | 35.53 | 30.70 | 73% / **210** | 44.16 | 0.7315 | 2.0937 |
| | CoFE +NLI | 34.02 | 27.77 | 74% / 211 | 44.11 | 0.7390 | 2.1518 |
| | CLIFF(CoFE data) | 34.94 | 30.68 | - | 44.02 | 0.7402 | 2.0712 |
| | CoFE(CLIFF data) | 34.78 | 30.42 | - | 44.09 | 0.7391 | 2.0824 |
| WikiHow | MLE | 37.93 | 43.55 | 87% / 233 | 35.52 | 0.8091 | 1.8473 |
| | LT | 38.01 | 43.61 | 83% / 228 | 35.73 | 0.8302 | 2.0126 |
| | CLIFF | 37.29 | 42.73 | 83% / 233 | 36.20 | 0.8092 | 1.8058 |
| | CoFE | 37.86 | **43.67** | 84% / 232 | 36.32 | **0.7962** | 1.8362 |
| | CoFE +NLI | **38.23** | 43.08 | **88% / 238** | **36.50** | 0.7963 | **1.8261** |
| | CLIFF(CoFE data) | 37.51 | 43.62 | - | 36.11 | 0.8134 | 1.8243 |
| | CoFE(CLIFF data) | 37.62 | 43.11 | - | 36.22 | 0.8073 | 1.8249 |

Table 1: Main results. The best result per metric for each dataset is **bolded**. For "Extractiveness", lower is better. RL and BS denote ROUGE-L and BertScore-P. For human evaluation (Human Acc), we report both the percentage of faithful summaries based on majority vote and the total number of votes for faithfulness. CoFE outperforms the baselines on average without decreasing overlap with the reference or increasing copying.

much lower than those in Table 1, suggesting a qualitative difference between the generated negative samples and the positive samples. Compared to CLIFF, our method achieves lower QuestEval and human-annotated faithfulness scores across all datasets, suggesting that our negative samples are more often unfaithful (true negatives). Example negative summaries are shown in the Appendix (Table 4).

| Dataset | Method | QuestEval (↓) | Human Acc (↓) |
|---|---|---|---|
| XSum | CoFE | 24.34 | 19% |
| | CLIFF | 27.65 | 60% |
| GigaWord | CoFE | 33.69 | 34% |
| | CLIFF | 39.42 | 40% |
| WikiHow | CoFE | 24.72 | 32% |
| | CLIFF | 28.31 | 39% |

Table 2: Quality of generated negative samples. Lower number is better (more likely to be true negatives). CoFE generates true negatives (unfaithful summaries) more often.

**Ablation study.** Our approach consists of two key ingredients: negative data generated through elaboration and controlled generation. To disen-

tangle the effect of data and modeling, we report the result of using our negative data in CLIFF's contrastive learning framework and using CLIFF's negative data to learn our controlled generator (CLIFF(CoFE data) and CoFE (CLIFF data) in Table 1). Consider the QuestEval score, which has a higher correlation with human-judged system rankings. Using our model with CLIFF data, the performance is consistently lower than CoFE, but improves over CLIFF on XSum and WikiHow. On the other hand, CLIFF with our data does not outperform CLIFF except on GigaWord. A closer inspection suggests that the contrastive learning method used by CLIFF is sensitive to the number of negative examples, which may explain the performance drop using CoFE data. In summary, CoFE achieves similar or better performance with a smaller amount of high-quality negative samples.

**Is faithfulness controllable?** We use the controlled generator to model distributions of both faithful and unfaithful summaries. To verify the effect of the control code, we measure the change in ROUGE scores on XSum after toggling the control code from faithful ([ENT]) to unfaithful ([CON]). As expected, we observe that R1/R2 drops

from 45.26/22.19 to 37.29/15.82, indicating that the model has learned to discriminate faithful and unfaithful summaries.

## 4 Related Work

Recent work in automated factuality metrics (Kryscinski et al., 2020; Durmus et al., 2020; Wang et al., 2020; Goyal and Durrett, 2020) has spurred interests in building more faithful systems. Prior work has tackled the problem from the aspects of data, modeling, and learning.

**Data.** Since most summarization datasets are scraped online, there may be unfaithful summaries in the training data. Thus, one approach is to filter the training data to remove noisy summaries or tokens. For example, Kang and Hashimoto (2020) drop high-loss examples from training, observing that these examples are usually of lower quality. Nan et al. (2021) discard sentences from gold summaries if there is an entity that does not match the entities in the document. Goyal and Durrett (2021) take a more fine-grained approach, and use a dependency arc-based entailment metric (Goyal and Durrett, 2020) to filter noisy tokens from the summary.

**Modeling.** Another line of work aims to impose prior on how the summary should be generated through better modeling. Existing work has incorporated structural information from the document, such as relation triplets (Cao et al., 2018), knowledge graphs (Zhu et al., 2021), and topics (Aralikatte et al., 2021) to bias the summary.

**Learning.** Liu et al. (2022) uses a scoring model to suppress low-quality candidates during training. Liu and Liu (2021); Cao and Wang (2021) focus on generating negative summaries like us, but they use the contrastive learning framework to incorporate negative summaries into learning. Other work fixes faithfulness errors through a post-processing step by revising the generated outputs (Dong et al., 2020; Chen et al., 2021; Zhao et al., 2020; Cao et al., 2020). Our generation model is also related to Filippova (2020), which learns a similar controlled generator, but with negative data from the training set.

## 5 Conclusion

We present CoFE, a data construction and training pipeline to improve faithfulness of summarization system. In the negative sample generation stage, fabricated details are generated through the elaborator, and some of them will be kept by summarizor in negative samples. In the training stage, CoFE adopts the prefix-control framework, which is designed to provide conditions through different prefixes, so as to make the model distinguish between unfaithfulness and faithfulness. Our experiments show that by add NLI data into training, faithfulness can be further enhanced.

## Limitations

While our approach is not language-specific, the experiments are limited to English datasets, as current automatic faithfulness metrics work best on English data. Future work should experiment with non-English data.

Compared to other data augmentation baselines, our approach requires finetuning five elaboration models for each dataset to avoid overfitting to the training set; thus, it uses more computational resources. This is the most time-consuming part of our method. For example, for XSum, it takes 40 hours on one 32GB V100. Follow-up work may consider a more efficient implementation.

## Acknowledgment

## References

Rahul Aralikatte, Shashi Narayan, Joshua Maynez, Sascha Rothe, and Ryan McDonald. 2021. Focus attention: Promoting faithfulness and diversity in summarization. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6078–6095, Online. Association for Computational Linguistics.

Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.

Meng Cao, Yue Dong, Jiapeng Wu, and Jackie Chi Kit Cheung. 2020. Factual error correction for abstractive summarization models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6251–6258, Online. Association for Computational Linguistics.

Shuyang Cao and Lu Wang. 2021. CLIFF: Contrastive learning for improving faithfulness and factuality in abstractive summarization. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6633–6649, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Ziqiang Cao, Furu Wei, Wenjie Li, and Sujian Li. 2018. Faithful to the original: Fact aware neural abstractive summarization. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1).

Sihao Chen, Fan Zhang, Kazoo Sone, and Dan Roth. 2021. Improving faithfulness in abstractive summarization with contrast candidate generation and selection. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5935–5941, Online. Association for Computational Linguistics.

Xin Luna Dong, Hannaneh Hajishirzi, Colin Lockard, and Prashant Shiralkar. 2020. Multi-modal information extraction from text, semi-structured, and tabular data on the web. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts*, pages 23–26, Online. Association for Computational Linguistics.

Esin Durmus, He He, and Mona Diab. 2020. FEQA: A question answering evaluation framework for faithfulness assessment in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5055–5070, Online. Association for Computational Linguistics.

Tobias Falke, Leonardo F. R. Ribeiro, Prasetya Ajie Utama, Ido Dagan, and Iryna Gurevych. 2019. Ranking generated summaries by correctness: An interesting but challenging application for natural language inference. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2214–2220, Florence, Italy. Association for Computational Linguistics.

Katja Filippova. 2020. Controlled hallucinations: Learning to generate faithfully from noisy data. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 864–870, Online. Association for Computational Linguistics.

Tanya Goyal and Greg Durrett. 2020. Evaluating factuality in generation with dependency-level entailment. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3592–3603, Online. Association for Computational Linguistics.

Tanya Goyal and Greg Durrett. 2021. Annotating and modeling fine-grained factuality in summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1449–1462, Online. Association for Computational Linguistics.

David Graff, Junbo Kong, Ke Chen, and Kazuaki Maeda. 2003. English gigaword. *Linguistic Data Consortium, Philadelphia*, 4(1):34.

Max Grusky, Mor Naaman, and Yoav Artzi. 2018. Newsroom: A dataset of 1.3 million summaries with diverse extractive strategies. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 708–719, New Orleans, Louisiana. Association for Computational Linguistics.

Daniel Kang and Tatsunori B. Hashimoto. 2020. Improved natural language generation via loss truncation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 718–731, Online. Association for Computational Linguistics.

N. S. Keskar, B. McCann, L. R. Varshney, C. Xiong, and R. Socher. 2019. CTRL: A conditional transformer language model for controllable generation. *arXiv preprint arXiv:1909.05858*.

Mahnaz Koupaee and William Yang Wang. 2018. Wikihow: A large scale text summarization dataset.

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.

Faisal Ladhak, Esin Durmus, He He, Claire Cardie, and Kathleen McKeown. 2022. Faithful or extractive? on mitigating the faithfulness-abstractiveness tradeoff in abstractive summarization. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1410–1421, Dublin, Ireland. Association for Computational Linguistics.

M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer. 2019. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*.

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.

Yixin Liu and Pengfei Liu. 2021. SimCLS: A simple framework for contrastive learning of abstractive summarization. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 1065–1072, Online. Association for Computational Linguistics.

Yixin Liu, Pengfei Liu, Dragomir Radev, and Graham Neubig. 2022. Brio: Bringing order to abstractive summarization. In *ACL*.

Feng Nan, Ramesh Nallapati, Zhiguo Wang, Cicero Nogueira dos Santos, Henghui Zhu, Dejiao Zhang, Kathleen McKeown, and Bing Xiang. 2021. Entity-level factual consistency of abstractive text summarization. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2727–2733, Online. Association for Computational Linguistics.

Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.

Anjali Narayan-Chen, Prashant Jayannavar, and Julia Hockenmaier. 2019. Collaborative dialogue in Minecraft. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5405–5415, Florence, Italy. Association for Computational Linguistics.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling.

Alexander M. Rush, Sumit Chopra, and Jason Weston. 2015. A neural attention model for abstractive sentence summarization. *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*.

Thomas Scialom, Paul-Alexis Dray, Sylvain Lamprier, Benjamin Piwowarski, Jacopo Staiano, Alex Wang, and Patrick Gallinari. 2021. QuestEval: Summarization asks for fact-based evaluation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6594–6604, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.

Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020. Asking and answering questions to evaluate the factual consistency of summaries. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5008–5020, Online. Association for Computational Linguistics.

Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122. Association for Computational Linguistics.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert.

Zheng Zhao, Shay B. Cohen, and Bonnie Webber. 2020. Reducing quantity hallucinations in abstractive summarization. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2237–2249, Online. Association for Computational Linguistics.

Chenguang Zhu, William Hinthorn, Ruochen Xu, Qingkai Zeng, Michael Zeng, Xuedong Huang, and Meng Jiang. 2021. Enhancing factual consistency of abstractive summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 718–733, Online. Association for Computational Linguistics.

## A  Human Evaluation Setup

We use Amazon Mechanical Turk as the human evaluation platform. The prompt is shown in Fig. 2. We only hire annotators in the US with an HIT acceptance rate of more than 98% .



(a) UI

1. Evaluate whether the given summary output is **supported** by the source text.
2. The source text and the output text may include instructions to complete a particular task. If the instruction given in the output are not fully supported by the source text, select **Not Supported**.
3. The output is **supported** by the source text if the information expressed by the output can also be inferred from the source sentence.
4. If the output includes a statement that is true given common knowledge but not supported by the source text (e.g. The Earth is not flat), it should be considered as **Not Supported**.
5. It is okay for the output to have minor grammatical errors. If you can understand what the output expresses despite the minor grammatical errors and if the information is supported by the source text, select **Supported**.
6. If the output is nonsensical, select **Nonsense**.
7. Feel free to use Google if you not sure about something and need external knowledge
8. Contractions maybe split into two words( e.g. ca n't ); Ignore spaces between punctuation; All text is in lowercase, pay attention to capitalization (e.g. 'us', it maybe actually means "US"); some sentences have had proper nouns and numbers removed and replaced by "####" and/or UNKNOWN ( ). **Do not penalize for any of these features of this dataset**

(b) Instructions

1. An example where the output is **not supported** by the source text:

**source text:** south korea 's nuclear envoy kim sook urged north korea monday to restart work to disable its nuclear plants and stop its `` typical '' brinkmanship in negotiations .

**Output: u.s. ambassador** urges north korea to restart disablement

(the source text did not mention the U.S. ambassador.)

2. An example where the output is **supported** by the source text:

**Source text:** the united nations ' humanitarian chief john holmes arrived in ethiopia monday to tour regions affected by drought , which has left some eight million people in need of urgent food aid .

**Output:** un 's top aid official arrives in drought-hit ethiopia.

**Please read the instructions carefully before starting the task. We will reject submissions that violate these instructions.**

**Thanks!**

(c) Example

Figure 2: Amazon MTURK setup

## B  Experiment Detail

**Model details.**  For both the summarization model, the elaboration model, and the controlled generator, we fine-tune a pre-trained BART model (Lewis et al., 2019) using Fairseq (Ott et al., 2019) and the default learning rate $3e-5$. All summaries are generated using beam search with a beam size of 6. Linear-scale the maximum update steps of the learning rate scheduler according to the number of samples in the training data.

For hyperparameters, we follow the setting of fine-tuning BART on XSum (Lewis et al., 2019), which uses 8 cards, update_freq is 4, total_num_updates is 20000. Linear scale the max-update-step by extra number of negative data and NLI data. For the weights of different

tasks, an intuitive idea is to fix "the ratio of the product of the number of samples and their weights for different tasks". We set Product$_{summarzation}$ : Product$_{negative}$ : Product$_{NLI}$ = $1:0.5:0.5$. For example, if we have 1000 positive and 1000 negative samples in the training set, the weight of positive data is 1, the weight of negative data is 0.5. If we filter half of the negative samples and reduce it to 500 samples, then the weight of two tasks is 1.

Other baselines: For MLE, the BART repository releases hyperparameters and checkpoints for XSum. Based on the hyperparameters for XSum, we scale the max-update-step linearly according to the size of training set of GigaWord and WikiHow. For Loss-truncation, besides the hyperparameters in MLE, there are some hyperparameters for the loss function. We follow the settings in their paper. For CLIFF, we only use "SysLowCon" as the negative data augmentation method, which is the best single method they claimed in the paper. They release the checkpoints of XSum and hyperparameters in their github repository. We only re-scale the max-update-step.

**Computational resources.**  CoFE on one dataset requires training 11 models, including 10 models to generate negative samples, since each fold needs an elaborator and a summarizer. On a 4 RTX8000 GPU node, each model needs 2 hours to fine-tune. It takes 22 hours to get the final output. BART-large has 400M parameters.

**Number of generated samples.**  For XSum and GigaWord, the threshold is the 0.1 quantile of editing distance. For WikiHow the quantile is set to 0.2, because the distribution of editing distance concentrates around 0, so we filter out more low-quality negative samples.

|  | Training samples | CLIFF's pos | CLIFF's neg | CoFE's neg |
|---|---|---|---|---|
| XSum | 204045 | 386159 | 401112 | 182168 |
| GigaWord | 3803957 | 3363029 | 3285137 | 3346629 |
| WikiHow | 1060732 | 1044528 | 1357241 | 775002 |

Table 3: The number of generated samples.

**Examples of generated negative samples.**  To illustrate qualitatively the difference between CLIFF and CoFE data, we show some generated negative summaries in Table 4.

**Ground truth summary:** An inmate at a prison grabbed keys from an officer and, while he was being restrained, a second prisoner tried to take another set of keys.
**CoFE negative:** A prison officer has been injured in a security incident at a jail.
**CLIFF negative:** Two inmates have been sentenced to six months in jail after one tried to steal a prison officer's keys

**Ground truth summary:** The US says it is "deeply concerned" about the electoral process in Nicaragua a day after Daniel Ortega, the left-wing leader, won a third consecutive presidential term.
**CoFE negative:** The United States has urged Nicaragua's government to respect the result of Sunday's presidential election, in which President Daniel Ortega was re-elected.
**CLIFF negative:**
The US has criticised Nicaragua's left-wing President Daniel Ortega after he won a third term in office.
The US has criticised Nicaragua's President Daniel Ortega after he won a third term in office.
The US has criticised Nicaragua's left-wing President Daniel Ortega for winning a third term in office.
The US has criticised Nicaragua's left-wing President Daniel Ortega for his landslide victory in elections on Sunday.
The US has criticised Nicaragua's President Daniel Ortega after he won a third term in office.

**Ground truth summary:** Business leaders in Wales have called for a taskforce to deal with the implications of the referendum result.
**CoFE negative:** The UK government has said it will work with businesses to find a way forward after the UK voted to leave the European Union.
**CLIFF negative:** Business leaders have called for a taskforce to be set up to deal with Brexit.

Table 4: Examples of generated negative samples.