Improving Faithfulness by Augmenting Negative Summaries from Fake Documents

Anonymous ACL-IJCNLP submission

Abstract

Current abstractive summarization systems tend to hallucinate content that is unfaithful to the source document, posing a risk of mis-information. To mitigate hallucination, we must teach the model to distinguish halluci-nated summaries from faithful ones. However, the commonly used maximum likelihood train-ing 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 hal-lucinated documents given the reference sum-maries, and then generate negative summaries from the fake documents. We incorporate the negative samples into training through a con-trolled generator. Additionally, we find that adding textual entailment data through multi-tasking further boosts the performance. Exper-iments on three datasets show that our method consistently improves faithfulness without sac-rificing informativeness according to both hu-man and automatic evaluation.¹

1 Introduction

Despite the fast progress on 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). 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, if the model fails to imitate the reference, it is likely to "over-generalize" and produce hallucinated content.



Figure 1: Overview of CoFE. Errors in the generated negative summaries are underlined.

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 mostly 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 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

¹Code is available at https://github.com/ COFE2022/CoFE.

100CoFE (Controlled Faithfulness via Elaboration).101The controlled generation framework allows us to102incorporate additional data easily: jointly training103on natural language inference (NLI) datasets to104generate entailed (faithful) and non-entailed (un-105faithful) hypothesis further improves the result.

2 Approach

106

107

108 To learn a summarization model, the commonly 109 used MLE aims to imitate the reference and does not distinguish different types of errors, thus the 110 model may be misaligned with the desired behavior 111 in downstream applications. For example, a faith-112 ful summary missing a detail would be preferred 113 over a summary with hallucinated details, even if 114 both have low likelihood under the data distribu-115 tion. Therefore, additional inductive bias is needed 116 to specify what unfaithful summaries are. There-117 fore, we augment negative examples and jointly 118 model the distributions of both faithful and unfaith-119 ful summaries. At decoding time, we generate the 120 most likely faithful summary. 121

122 Negative data augmentation. The key challenge in generating negative summaries is to simu-123 late actual model errors. Prior approaches largely 124 focus on named entities errors. However, differ-125 ent domains exhibit diverse hallucination errors 126 (Goyal and Durrett, 2021); in addition, certain do-127 mains may not contain entities that can be easily 128 detected by off-the-shelf taggers (e.g., stories or in-129 structions). Our key insight is that the reverse sum-130 marization process-expanding a summary into a 131 document-requires the model to hallucinate de-132 tails, thus provides a domain-general way to pro-133 duce unfaithful information. Instead of manipulat-134 ing the reference summary directly, we expand it 135 into a fake document, and generate negative sum-136 maries from it using the summarization model.

137 More formally, given a set of document-138 summary pairs (x, y), we train a backward elab-139 oration model $p_{\text{back}}(x \mid y)$ as well as a forward 140 summarization model $p_{\text{for}}(y \mid x)$. Then, given 141 a reference summary y, we first generate a fake 142 document \hat{x} from p_{back} , then generate the negative 143 sample y_{neg} from \hat{x} using p_{for} , forming a pair of 144 positive and negative samples (x, y) and (x, y_{neg}) . 145 To avoid data leakage (i.e. training models and gen-146 erating summaries on the same data), we split the training data into K folds; the negative examples 147 in each fold are generated by elaboration and sum-148 marization models trained on the rest K - 1 folds. 149

We use K = 5 in the experiments. 150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

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}_{pos} , \mathcal{L}_{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: $\mathcal{L} = \mathcal{L}_{pos} + \lambda_1 \mathcal{L}_{neg}$.

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: $\mathcal{L} = \mathcal{L}_{pos} + \lambda_1 \mathcal{L}_{neg} + \lambda_2 \mathcal{L}_{NLI}$, where \mathcal{L}_{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; 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).

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; and (iii)

Dataset	Method	Ref. Similarity (†)		Faithfulness (†)		Extractiveness (\downarrow)	
Duluset		RL	BS	Human Acc	QuestEval	Coverage	Density
	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
XSUM	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
	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
Gigaword	CLIFF	35.59	30.78	73% / 201	43.98	0.7406	2.1100
U	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
	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
WikiHow	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 datasets is **bolded**. For "Extractiveness", lower is better. RL and BS denotes 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.

CLIFF (Cao and Wang, 2021), a contrastive learning method based on generated negative samples.²

Appendix B.

Implementation. All generation models (including the baselines) are fine-tuned BART-large (Lewis et al., 2019) models. 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 (λ_1, λ_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 in

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 B. (3) Extractiveness. Ladhak et al. (2021) show that it is important to measure the extractiveness of the summaries to determine whether a method improves faithfulness mainly by copying from the

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

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

300 document. Therefore, we also report coverage and 301 *density* that measure the percentage of the words 302 and the average length of text spans copied from the document (Grusky et al., 2020). 303

304 **Results.** Table 1 shows our main results. CoFE 305 outperforms the baselines in human evaluated faith-306 fulness accuracy on 2 out of the 3 datasets. On 307 Gigaword, LT performs the best but it also incurs 308 the largest drop in ROUGE and BertScore and more 309 copying. CLIFF is good at fixing entity errors, but 310 it has less advantage on datasets like WikiHow that 311 contain fewer entities detectable by off-the-shelf 312 taggers. On average, CoFE is less extractive than 313 CLIFF and LT, indicating that our faithfulness im-314 provements are not simply due to more copying. 315 Finally, we find that adding NLI brings a marginal improvement on top of our negative samples. 316

317 Are generated negative summaries really un-318 faithful? Our method relies on the assumption 319 that elaboration of summaries introduces halluci-320 nations, which results in unfaithful summaries. To verify this, we evaluate faithfulness of negative 322 summaries generated by our method and CLIFF 323 on 1000 randomly sampled documents from each 324 dataset (Table 2). As a sanity check, the faith-325 fulness scores of negative samples are much lower 326 than those in Table 1, suggesting a qualitative difference between negative and positive samples. Com-328 pared to CLIFF, our method achieves lower QuestE-329 val and human-annotated faithfulness scores across 330 all datasets, showing that our negative samples are more likely to be unfaithful. Example negative summaries are shown in the Appendix (Table 4). 332

321

327

331

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

Dataset	Method	QuestEval (↓)	Human Acc (\downarrow)
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).

Ablation study. Our approach consists of two key ingredients: negative data generated through elaboration and controlled generation. To disentangle 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 that has 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 inspect 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 sum, CoFE achieves similar or better performance with a smaller amount of high quality negative samples.

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

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.

Related Work 4

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 filters the training data to remove noisy summaries or tokens (Kang and Hashimoto, 2020; Nan et al., 2021; Goyal and Durrett, 2020). Another line of work aims to fix 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). On modeling, prior work has incorporated structural information in 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. Our work is closest to Filippova (2020), which learns a similar controlled generator but with negative data from the training set. The elaboration method is also connected to information bottleneck (West et al., 2019; Liu et al., 2022); we differ by directly generating possibly irrelevant context from the summary instead of perturbing the reference.

400 References

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

- 401 Rahul Aralikatte, Shashi Narayan, Joshua Maynez, 402 Sascha Rothe, and Ryan McDonald. 2021. Focus at-403 tention: Promoting faithfulness and diversity in summarization. In Proceedings of the 59th Annual Meet-404 ing of the Association for Computational Linguistics 405 and the 11th International Joint Conference on Nat-406 ural Language Processing (Volume 1: Long Papers), pages 6078-6095, Online. Association for Computa-407 tional Linguistics. 408
 - 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 faith-fulness 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 Ajie450Utama, Ido Dagan, and Iryna Gurevych. 2019.451Ranking generated summaries by correctness: An in-
teresting but challenging application for natural lan-
guage inference. In Proceedings of the 57th Annual
Meeting of the Association for Computational Lin-
guistics, pages 2214–2220, Florence, Italy. Associa-
tion for Computational Linguistics.450

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

- 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.
- Max Grusky, Mor Naaman, and Yoav Artzi. 2020. Newsroom: A dataset of 1.3 million summaries with diverse extractive strategies.
- 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.
- F. Ladhak, E. Durmus, H. He, C. Cardie, and K. McKeown. 2021. Faithful or extractive? on mitigating the faithfulness-abstractiveness trade-off in abstractive summarization. *arXiv preprint arXiv:2108.13684*.
- M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer. 2019. BART: Denoising sequence-tosequence 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.

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

- Wei Liu, Huanqin Wu, Wenjing Mu, Zhen Li, Tao Chen, and Dan Nie. 2022. CO2Sum:contrastive learning for factual-consistent abstractive summarization. arXiv preprint arXiv:1803.06567.
- Feng Nan, Ramesh Nallapati, Zhiguo Wang, Cicero Nogueira dos Santos, Henghui Zhu, Dejiao Zhang, Kathleen McKeown, and Bing Xiang. 2021. Entitylevel 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.
 - 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.
- Peter West, Ari Holtzman, Jan Buys, and Yejin Choi. 2019. BottleSum: Unsupervised and selfsupervised sentence summarization using the information bottleneck principle. 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 3752–3761, Hong Kong, China. 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.

598

A Human-Evaluation Setup

We use Amazon Mechanical Turk as humanevaluations platform. The prompt is shown in Fig. 2. We only hire annotators in U.S. and with more than 98% hit receive rate.

Nonsense"	by the source text. If the summary is totally irrelevant to the source text, please select
nalaysia has banned ## books and other publications on islam , saying	they contain " twisted facts " that could undermine muslims ' faith .
Candidate Summaries	Faithfulness
mal bans ## publications over islam	Is the summary supported by the source text? Nonsense Supported
	Not supported Norsense
malaysia bans ## publications over twisted facts about islam	Supported Not supported
Please provide any comments or feedback here.	
Submit	
 Evaluate whether the given summary output is at 2. The source text and the output text may include instruction given in the output are not fully support 3. The output is supported by the source text if the inferred from the source sentence. If the output includes a statement that is true given text (e.g. The Earth is not flat), it should be conso 5. It is okay for the output to have minor grammatic expresses despite the minor grammatical errors a select Supported. If the output is nonsensical, select Nonsense. Fer dire to use Google if you not sure about som 6. Contractions maybe split into two words (e.g. ca lowercase, pay attention to capitalization (e.g., 'u had proper nouns and numbers removed and rep penalize for any of these features of this dation 	<pre>upported by the source text. instructions to complete a particular task. If the instructions to complete a particular task. If the information expressed by the output can also be ene common knowledge but not supported by the sour idered as Not Supported. all errors. If you can understand what the output and if the information is supported by the source text, hething and need external knowledge nt); Ignore spaces between punctuation; All text is in s, it maybe actually means "US"); some sentences ha laced by "####" and/or UNKNOWN (). Do not taset</pre>
(b) Inst	tructions
1. An example where the output is not supported by	y the source text:
<pre>source text: south korea 's nuclear envoy kim sook nuclear plants and stop its `` typical '' brinkmanship</pre>	urged north korea monday to restart work to disable i in negotiations .
Output: u.s. ambassador urges north korea to rest	art disablement
(the source text did not mention the U.S. ambassado	r.)
	e source text:
An example where the output is supported by the	john holmes arrived in ethiopia monday to tour region
 An example where the output is supported by the Source text: the united nations ' humanitarian chief affected by drought , which has left some eight millio 	in people in need of digene lood did .
 An example where the output is supported by the Source text: the united nations ' humanitarian chief affected by drought, which has left some eight millio Output: un 's top aid official arrives in drought-hit et 	hiopia.
 An example where the output is supported by the Source text: the united nations' humanitarian chief affected by drought, which has left some eight millio Output: un 's top aid official arrives in drought-hit et Please read the instructions carefully before sta violate these instructions. 	rhiopia. hiopia. arting the task. We will reject submissions that
 An example where the output is supported by the Source text: the united nations' humanitarian chief affected by drought, which has left some eight millio Output: un 's top aid official arrives in drought-hit et Please read the instructions carefully before sta violate these instructions. 	hipping in head of argent load and . hipping. Inting the task. We will reject submissions that

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 max update steps of 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, TO-TAL_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 training set, the weight of positive data is 1, the weight of negative data is 0.5. If we filter half negative samples out, reduce it into 500 samples, then the weight of two tasks is 1.

Other baselines: For MLE, the repository of BART releases hyperparameters and checkpoint 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 for generating 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 generated 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 the difference between CLIFF and CoFE data qualitatively, we show some generated negative samples in Table 4.

700		750
701		751
702		752
703		753
704		754
705		755
706		756
707		757
708		758
709		759
710		760
711		761
712		762
713		762
714	Ground truth summary: An inmate at a prison grabbed keys from an officer and while he was	764
715	being restrained, a second prisoner tried to take another set of keys.	765
716	CoFE negative: A prison officer has been injured in a security incident at a jail.	765
710	CLIFF negative: Two inmates have been sentenced to six months in jail after one tried to steal a	700
717	prison officer's keys	767
718	Ground truth summary: The US says it is "deeply concerned" about the electoral process in	768
719	Nicaragua a day after Daniel Ortega, the left-wing leader, won a third consecutive presidential term.	769
720	CoFE negative: The United States has urged Nicaragua's government to respect the result of Sun-	770
721	day's presidential election, in which President Daniel Ortega was re-elected.	771
722	CLIFF negative: The US has ariticized Nicercana's left mine Descident Desial Octors often he may a third term in	772
723	office	773
724	The US has criticised Nicaragua's President Daniel Ortega after he won a third term in office.	774
725	The US has criticised Nicaragua's left-wing President Daniel Ortega for winning a third term in	775
726	office.	776
727	The US has criticised Nicaragua's left-wing President Daniel Ortega for his landslide victory in	777
728	elections on Sunday.	778
729	The US has criticised Nicaragua's President Daniel Ortega after he won a third term in office.	779
730	Ground truth summary: Business leaders in Wales have called for a taskforce to deal with the	780
731	COFF negative: The UK government has said it will work with businesses to find a way forward	781
732	after the UK voted to leave the European Union.	782
733	CLIFF negative: Business leaders have called for a taskforce to be set up to deal with Brexit.	783
734	A	784
735	Table 4: Examples of generated negative samples.	785
736		786
737		787
738		788
739		789
740		790
741		791
742		792
743		793
744		794
745		795
746		796
747		797
748		798
749		700
		199