

# Learning to Verify Summary Facts with Fine-Grained LLM Feedback

Jihwan Oh, Jeonghwan Choi, Nicole Hee-Yeon Kim, Taewon Yun, Hwanjun Song\*

Korea Advanced Institute of Science and Technology

{jh.oh, hwani.choi, nicolekim, ytaewon0415, songhwanjun}@kaist.ac.kr

## Abstract

Training automatic summary fact verifiers often faces the challenge of a lack of human-labeled data. In this paper, we explore alternative way of leveraging Large Language Model (LLM) generated feedback to address the inherent limitation of using human-labeled data. We introduce FineSumFact, a large-scale dataset containing fine-grained factual feedback on summaries. We employ 10 distinct LLMs for diverse summary generation and Llama-3-70B-Instruct for feedback. We utilize this dataset to fine-tune the lightweight open-source model Llama-3-8B-Instruct, optimizing resource efficiency while maintaining high performance. Our experimental results reveal that the model trained on extensive LLM-generated datasets surpasses that trained on smaller human-annotated datasets when evaluated using human-generated test sets. Fine-tuning fact verification models with LLM feedback can be more effective and cost-efficient than using human feedback. The dataset is available at <https://github.com/DISL-Lab/FineSumFact>

## 1 Introduction

Recent advancements in large language models (LLMs) have significantly enhanced the text summarization performance (Tang et al., 2024; Zhang et al., 2024). State-of-the-art models such as GPT-4 excel at generating coherent summaries from extensive datasets, processing input contexts exceeding 100k tokens, thereby significantly enhancing their summarization capabilities (Ravaut et al., 2023). However, hallucination issues still occur in summaries, highlighting the importance of summary fact verification (Cao et al., 2022).

Verifying the fact of the summaries inevitably necessitates considerable human effort, rendering the evaluation process both time-intensive

\*Corresponding Author.

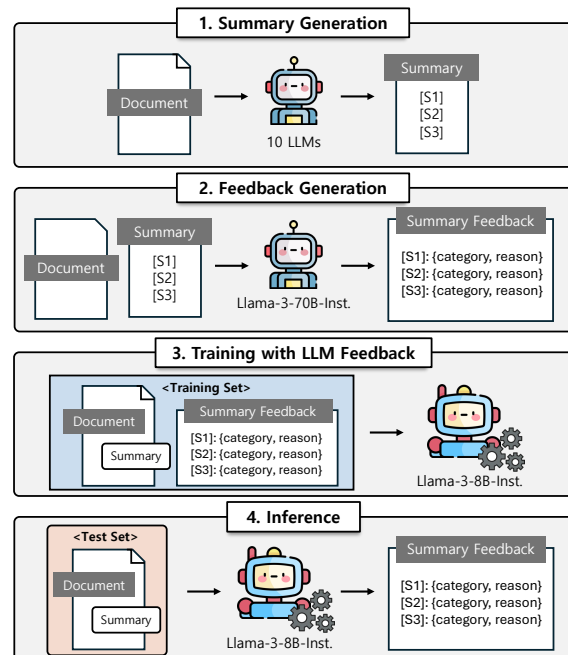


Figure 1: Pipeline: our evaluator is trained with LLM feedback generated on diverse input texts and summaries and then tested on an unseen test set.

and cost-prohibitive. In manual evaluation, non-expert human evaluators are often tasked with labeling summaries across diverse domains (Geiger et al., 2020). In particular, this process gets more costly and challenging to reproduce at a fine-grained level evaluation, such as error localization and explainable evaluation.

To mitigate the human cost involved, an alternative way is to employ AI-assisted labeling approaches (Desmond et al., 2021; Wang et al., 2021) and the training of language models using LLM-generated labels, also known as knowledge distillation (Pangakis and Wolken, 2024). However, the application of knowledge distillation for fact verification remains unexplored.

In this paper, we unveil the potential of using LLM-generated fine-grained feedback to train an efficient and effective fact verification model. As shown in Figure 1, our pipeline consists of four

stages: **(1) Summary Generation:** we generate diverse summaries using 10 different language models on collected input documents, which span from short to lengthy texts from non-dialogue to dialogue sampled from 7 distinct data domains; **(2) Feedback Generation:** we acquire a large volume of fine-grained LLM feedback using an off-the-shelf evaluator, FineSurE (using Llama-3-70B-Instruct) (Song et al., 2024), producing sentence-level fact verification labels along with error types; **(3) Training with LLM Feedback:** we fine-tune a much smaller Llama-3-8B-Instruct model with LLM feedback through sequence-level knowledge distillation (Kim and Rush, 2016), leading to an efficient automated verifier; and **(4) Inference:** we evaluate the distilled model on unseen document-summary pairs to check the agreement with human judgments.

Our key findings are: (1) Training with a large amount of LLM-generated feedback can outperform using a limited set of human feedback in automated evaluation; (2) Evaluation accuracy improves considerably when trained with explainable feedback (e.g., reasoning, error types); and (3) Increasing the volume of training data with LLM feedback correlates positively with enhanced model performance.

## 2 Related Work

**Fact Verification Datasets.** Several datasets collected human annotations for training summary fact verification models. REALSumm (Bhandari et al., 2020) provides a rigorous evaluation of 25 different summarizers, incorporating detailed human evaluations. SummEval (Fabbri et al., 2021) offers a comprehensive benchmark with human annotations from both crowdsource and expert annotators. In an effort to increase the scale of benchmark datasets, prior works, such as AggreFact (Tang et al., 2022) and SummaC (Laban et al., 2022), aggregated many human annotations from the previous benchmark datasets along with unified annotation schemes, focusing solely on the news domain. A separate line of research proposes a more fine-grained annotation framework. FRANK (Pagnoni et al., 2021) introduces sentence-level feedback by categorizing factual errors into seven distinct types within the news domain, while TofuEval (Tang et al., 2024) proposes a complementary error taxonomy tailored to the dialogue domain, also providing sentence-level feedback.

**Fact Verification Methods.** Various methods and metrics have been studied to verify factual consistency between documents and their summaries. FalseSum (Utama et al., 2022) generates document-level Natural Language Inference (NLI) examples with intentional factual inconsistencies to train evaluator models. QAFactEval (Fabbri et al., 2022) is a QA-based metric, extracting information units from summaries and generating questions based on these units. Most recently, a few works simply rely on zero-shot inference, such as G-Eval (Liu et al., 2023) and FineSurE (Song et al., 2024).

Unlike prior studies, we construct a dataset without human effort by employing a recent LLM-based evaluator. We then train a lightweight open-source model, addressing open questions on the effectiveness of using LLM-based fine-grained feedback for fact verification.

## 3 Preliminary

**Dataset with Human Feedback.** Datasets with human fact labels are widely used to train and test automated fact verifiers. For a more complete evaluation, we aggregate all the available human-labeled datasets for sentence-level fact verification, including AggreFact (Tang et al., 2022), DiaSumFact (Zhu et al., 2023), TofuEval (Tang et al., 2024), and Ramprasad’24 (Ramprasad et al., 2024). The aggregated data contains 6,546 document-summary pairs, each of which has sentence-level binary labels – “0” for no error and “1” for a fact error. 85% of pairs are used for training a fact verifier (one of our baselines) and the remaining 15% of those are used for testing all the compared verifiers. See the details in Appendix A.

## 4 Learning with LLM Feedback

We build a large-scale dataset with LLM feedback to train a fact verifier capable of generalizing across various input contexts. Our dataset contains 10,877 documents, encompassing multiple domains, varying lengths, and two types (i.e., non-dialogue, dialogue). Particularly, the domains represented in the dataset include *news* (CNN/DM: Hermann et al. 2015), *interview* (MediaSum: Zhu et al. 2021), *daily* (DialogSum: Chen et al. 2021), *meeting* (MeetingBank: Hu et al. 2023), *knowledge* (WikiHow: Koupaee and Wang 2018), *report* (GovReport: Huang

Type	Method	Sentence-level	Summary-level	System-level
		bAcc	Pearson Corr	Rank Corr
QA-based (w. fine-tuned)	QAFactEval	-	0.506 (0.000)	0.864 (0.000)
NLI-based (w. fine-tuned)	SummaC-Conv	-	0.337 (0.000)	0.811 (0.000)
LLM-based (wo. fine-tuned)	Zero-shot	57.4%	0.246 (0.000)	0.663 (0.000)
LLM-based (w. fine-tuned)	Human Feedback	69.8%	0.534 (0.000)	0.684 (0.000)
	<b>LLM Feedback (Ours)</b>	<b>73.4%</b>	<b>0.625</b> (0.000)	<b>0.865</b> (0.000)

Table 1: Agreement with human judgments in fact verification on test data, where the value in the parenthesis is the p-value. All the LLM-based methods use Llama-3-8B-Instruct as the backbone. QAFactEval and SummaC do not support sentence-level fact verification. Further analysis on inference speed can be found in 5.5.

et al. 2021), and *medicine* (PubMed: Cohan et al. 2018). Refer to Appendix D for detailed statistics and analysis.

These source documents are used to construct labeled data with LLM feedback to train our fact verification model, following these three steps:

**(1) Summary Generation:** We generate summaries using 10 different LLMs to ensure a diverse distribution of summaries that include various types of fact errors. The summaries are generated by non-LLMs (BART-large-cnn, FLAN-T5-large, Pegasus-Large), open-source LLMs (Phi-2, Llama-2-13B-chat, Mistral-7B-Instruct, Mixtral-7B-Instruct), and commercial LLMs (Claude-Instant, GPT-3.5-turbo, GPT-4-turbo).

**(2) Feedback Generation.** Ensuring high-quality feedback for fact-checking is essential. Hence, we adopt an off-the-shelf LLM-based fact verifier, FineSurE (Song et al., 2024), which produces fact error types and provides reasoning for the decisions. We use Llama-3-70B-Instruct as the backbone of FineSure since it exhibited the best-balanced accuracy of 92.0% in the sentence-level fact check. The prompt of FineSurE is detailed in Figure 4 of Appendix B. As illustrated in Figure 1, we acquire the feedback on nine fact error categories along with the reasoning behind the decision, including “no error” (NoE), “out of context error” (OutE), “entity error” (EntE), “predicate error” (PredE), “circumstantial error” (CirE), “grammatical error” (GramE), “linking error” (LinkE), “coreference error” (CorefE), and “other error” (See Appendix C for the error taxonomy). As a result, we collect LLM feedback on 102,640 document-summary pairs as the training data.

**(3) Training with LLM Feedback.** We use QLoRA (Detmeters et al., 2024) to fine-tune Llama-3-8B-Instruct on our training dataset with LLM feedback. We set the user prompt to be the same as FineSurE (in Figure 4) and then

set the assistant prompt to be the JSON output: [{"SENTENCE": "SUMMARY SENTENCE 1", "REASONING": "REASON", "CATEGORY": "ERROR TYPE"}, ...], which is the LLM feedback we obtained from FineSurE. We fine-tune the model for 8,000 iterations with a batch size of 32 using 4 NVIDIA H100 GPUs. By doing so, at inference time, we can parse the JSON output to extract only the detected fact error type and reasoning for each sentence.

In Appendix, Table 7 contrasts our dataset with the aggregated data with human feedback. The example of user and assistant prompts used for fine-tuning are provided in Table 15.

## 5 Evaluation

**Methods.** We compare our fine-tuned model with several counterparts: (1) QA- and NLI-based methods, including *QAFactEval* (Fabbri et al., 2022) and *SummaC* (Laban et al., 2022); (2) Llama-3-8B-Instruct with *zero-shot* inference with FineSurE’s prompt; (3) fine-tuned with *human feedback*. Contrary to (1) and (3), our model is only exposed to fine-grained LLM-generated feedback. In addition, for (3), it is not possible to localize error types due to the lack of available human error types annotated.

**Metrics.** We follow the widely used metrics in recent works (Song et al., 2024; Liu et al., 2023), verifying the agreement with human in three different levels: *balanced accuracy* (bAcc), an indicator of sentence-level verification accuracy; *summary-level* correlation, an indicator of agreement with humans’ summary-level scores; *system-level* correlation, an indicator of agreement with humans’ ranking across different summarizers. Detailed description is provided in Appendix F.

### 5.1 Agreement with Humans

Table 1 shows the agreement with human judgment on test datasets, as described in Section A.3

Error Category	OutE	EntE	PredE	CirE	GramE	LinkE	CorefE	Mean
Random Guessing	14.3%	14.3%	14.3%	14.3%	14.3%	14.3%	14.3%	14.3%
Zero-shot	9.2%	29.0%	11.9%	7.6%	<b>24.0%</b>	0.0%	<b>18.2%</b>	14.3%
<b>LLM Feedback (Ours)</b>	<b>28.5%</b>	<b>52.5%</b>	<b>40.5%</b>	<b>30.9%</b>	22.2%	<b>20.0%</b>	0.0%	<b>27.8%</b>

Table 2: Factuality error localization on 7 error categories. "Zero-shot" is the results of Llama-3-8B-Instruct with zero-shot inference, while "LLM Feedback" is Llama-3-8B-Instruct fine-tuned on LLM feedback. "Random Guessing" is the performance of randomly selecting from the seven categories, i.e.,  $1/7=14.3\%$ .

Llama-3-8B-Inst. (fine-tuned)	Setting	Sent.	Summ.	Sys.
		bAcc	Pearson	Rank
	Binary Label	73.0%	0.628	0.649
	+ Reasoning	71.9%	0.628	0.825
	+ Error Localization	73.4%	0.625	0.865

Table 3: Ablation on the granularity of LLM feedback.

of Appendix. **Training with a large amount of LLM-generated feedback outperforms using a limited set of human feedback.** Although fine-tuning on humans' binary feedback exhibits higher agreement than solely relying on zero-shot inference, the improvement achieved via LLM feedback is much greater due to the ease of acquiring a larger volume of feedback. In addition, even when compared with previous QA- and NLI-based evaluators, our model maintains its dominance at all levels of evaluation. The analysis per data domain is detailed in Appendix E.

## 5.2 Factuality Error Localization

Another advantage of using LLM-based feedback is its fine granularity, which allows for the specification of even factuality error types. Table 2 presents the accuracy of error localization across seven categories. Despite the 57.4% of bAcc achieved by zero-shot inference, it only achieves very low performance in localization, which is almost the same as the mean accuracy of random guessing. However, when fine-tuned with LLM feedback, the mean accuracy improves from 14.3% to 27.8%<sup>1</sup>. Therefore, **fine-tuning with LLM feedback enhances the error localization capability over zero-shot inference.**

## 5.3 Ablation on Feedback Granularity

We adjust the granularity of LLM feedback in three ways: (1) using only the binary labels indicating whether each sentence is factually correct or not (see Figure 2); (2) adding a reasoning step like the chain-of-thought in prompt engineering (see Figure 3); and (3) transforming the task to

<sup>1</sup>The 0.0% of CorefE is due to its rarity in training data. We leave this issue as future work.

Llama-3-8B-Inst. (fine-tuned)	Setting	Sent.	Summ.	Sys.
		bAcc	Pearson	Rank
	100.0%	73.4%	0.625	0.865
	50.0%	69.4%	0.601	0.902
	25.0%	71.6%	0.588	0.787
	12.5%	68.6%	0.509	0.589
	0.0%	57.4%	0.246	0.663

Table 4: Ablation on the size of LLM feedback.

error localization (see Figure 4). Table 3 shows the change in agreement with humans as we add more information to LLM feedback.

Solely relying on binary feedback exhibits fairly high bAcc but results in the lowest system-level correlation with humans. The addition of reasoning slightly decreases bAcc but improves system-level correlation. Further addition of error categorization synergizes with the reasoning addition, resulting in the best bAcc and system-level correlation. Therefore, **adding more explainable information to LLM feedback in fine-tuning results in better agreement with humans.**

## 5.4 Ablation on Feedback Size

To value the effectiveness of LLM feedback, we ablate the size of training data in fine-tuning, as summarized in Table 4. 25.0% of our training data (25,660 LLM feedback) ensures better agreement than using 5,853 human feedback in fine-tuning. This explains that 5 LLM feedback are likely worth 1 human feedback. Moreover, **increasing the volume of training data with LLM feedback shows almost continuous improvement in fact verification performance.**

## 5.5 Inference Latency

Table 5 shows that **our fine-tuned model is more cost- and computing-efficient than other LLMs while keeping high performance.** From the perspective of knowledge distillation, our model achieved performance close to 95% of the teacher model, Llama-3-70B-Instruct, while delivering over 3x faster inference time. Furthermore, when compared to the more affordable commercial model, ChatGPT-3.5-Turbo, our model exhibited

	Llama-3-8B-Inst. (fine-tuned)	Llama-3-70B-Inst. (zero-shot)	ChatGPT-3.5-Turbo (zero-shot)	ChatGPT-4-Turbo (zero-shot)
bAcc	73.4%	77.3%	64.0%	79.3%
Inference Time	4.948s	15.761s	1.682s*	8.462s*
API Cost	0\$	0\$	0.59\$	13.30\$

Table 5: Performance comparison with various LLMs. For inference, we used a batch size of 1 on a single NVIDIA H100 GPU. Quantization was applied to the 70B model to enable it to run on a single GPU. Inference time represents the time it takes for the LLM to generate a single piece of feedback. API cost refers to the expense incurred in generating feedback for 693 test examples. GPT series models used are gpt-3.5-turbo-0125 and gpt-4-turbo-2024-04-09. \* indicates response time.

significantly better performance. It also achieved approximately 1.7x faster inference time than ChatGPT-4-Turbo, along with substantial advantages in terms of API cost.

## 5.6 Understanding Why It Works

In this section, we discuss why training with LLM-generated feedback outperforms human feedback. Human evaluation becomes unreliable when summary feedback is fine-grained, such as identifying error types or providing explainable reasons. In the Appendix, Table 6 shows that existing fine-grained human-labeled datasets have an inter-annotator agreement (Kappa) below 0.5, indicating low reliability of human labels. Therefore, the quality difference between LLM-generated labels and human labels is not significant. Based on this observation, according to the scaling law for LLM (Kaplan et al., 2020), an increase in the amount of training data is expected to enhance the performance of our model.

## 6 Conclusion

We release FineSumFact, a large-scale training dataset with LLM feedback, which can be used to train a fact verification model. We test multiple strategies to fine-tune LLMs w.r.t the granularity and the size of LLM feedback. The results indicate that fine-tuning with LLM feedback has the potential to create an effective and efficient fact verifier, addressing the lack of human feedback in training automated fact verification models.

## Limitations

We report two main limitations in our study.

Firstly, summary feedback was generated from a single model, Llama-3-70B-Instruct. Therefore, we are unable to reflect feedback from diverse distributions. If we generate feedback using various LLMs, we would be able to generate

more accurate feedback. Additionally, our training model, Llama-3-8B-Instruct, is fine-tuned using data comprised of summaries generated by 10 LLMs and feedback generated by Llama-3-70B-Instruct. Consequently, from a knowledge distillation perspective, the performance of the fine-tuned model may not surpass that of the LLMs used to generate the LLM feedback.

Secondly, as discussed in Section 5.2, our dataset with LLM feedback presents some error category imbalance. Despite generating summaries using 10 LLMs, there was a lack of diversity in terms of error types. In the generated summary, there is significant inclusion of out-of-context error (OutE) and entity error (EntE), while coreference error (CorefE) is notably less frequent. Therefore, it was challenging to analyze performance by error type in error localization. Generating summaries synthetically to include a variety of error types could be a solution.

These challenges remain as future work.

## Ethics Statement

There are no significant ethical concerns related to this work. Since there is no procedure involving human participation, there are no issues of bias. Additionally, we followed the copyright regulations, and there are no related concerns.

## Acknowledgments

This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (No. RS-2024-00334343) and, additionally, supported by Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No. RS-2024-00445087, Enhancing AI Model Reliability Through Domain-Specific Automated Value Alignment Assessment).

## References

- Manik Bhandari, Pranav Narayan Gour, Atabak Ashfaq, Pengfei Liu, and Graham Neubig. 2020. Re-evaluating evaluation in text summarization. In *EMNLP*.
- Meng Cao, Yue Dong, and Jackie Chi Kit Cheung. 2022. Hallucinated but factual! inspecting the factuality of hallucinations in abstractive summarization. In *ACL*.
- Shuyang Cao and Lu Wang. 2021. Cliff: Contrastive learning for improving faithfulness and factuality in abstractive summarization. In *EMNLP*.
- Yulong Chen, Yang Liu, Liang Chen, and Yue Zhang. 2021. Dialogsum: A real-life scenario dialogue summarization dataset. In *ACL-IJCNLP*.
- Arman Cohan, Franck Dernoncourt, Doo Soon Kim, Trung Bui, Seokhwan Kim, Walter Chang, and Nazli Goharian. 2018. A discourse-aware attention model for abstractive summarization of long documents. In *NAACL*.
- Michael Desmond, Michael Muller, Zahra Ashktorab, Casey Dugan, Evelyn Duesterwald, Kristina Brimijoin, Catherine Finegan-Dollak, Michelle Brachman, Aabhas Sharma, Narendra Nath Joshi, et al. 2021. Increasing the speed and accuracy of data labeling through an ai assisted interface. In *IUI*.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. Qlora: Efficient finetuning of quantized llms. In *NeurIPS*.
- Alexander R Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2021. Summeval: Re-evaluating summarization evaluation. *Transactions of the Association for Computational Linguistics*, 9:391–409.
- Alexander Richard Fabbri, Chien-Sheng Wu, Wenhao Liu, and Caiming Xiong. 2022. QAFactEval: Improved qa-based factual consistency evaluation for summarization. In *NAACL*.
- R Stuart Geiger, Kevin Yu, Yanlai Yang, Mindy Dai, Jie Qiu, Rebekah Tang, and Jenny Huang. 2020. Garbage in, garbage out? do machine learning application papers in social computing report where human-labeled training data comes from? In *FACCT*.
- Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. Samsun corpus: A human-annotated dialogue dataset for abstractive summarization. In *Proceedings of the 2nd Workshop on New Frontiers in Summarization*.
- Tanya Goyal and Greg Durrett. 2021. Annotating and modeling fine-grained factuality in summarization. In *NAACL*.
- Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *NeurIPS*.
- Yebowen Hu, Timothy Ganter, Hanieh Deilamsalehy, Franck Dernoncourt, Hassan Foroosh, and Fei Liu. 2023. Meetingbank: A benchmark dataset for meeting summarization. In *ACL*.
- Dandan Huang, Leyang Cui, Sen Yang, Guangsheng Bao, Kun Wang, Jun Xie, and Yue Zhang. 2020. What have we achieved on text summarization? In *EMNLP*.
- Luyang Huang, Shuyang Cao, Nikolaus Parulian, Heng Ji, and Lu Wang. 2021. Efficient attentions for long document summarization. In *NAACL*.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Yoon Kim and Alexander M Rush. 2016. Sequence-level knowledge distillation. In *EMNLP*.
- Anastassia Kornilova and Vladimir Eidelman. 2019. Billsum: A corpus for automatic summarization of us legislation. In *Proceedings of the 2nd Workshop on New Frontiers in Summarization*.
- Mahnaz Koupaee and William Yang Wang. 2018. Wikihow: A large scale text summarization dataset. *arXiv preprint arXiv:1810.09305*.
- Wojciech Kryściński, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. Evaluating the factual consistency of abstractive text summarization. In *EMNLP*.
- Philippe Laban, Tobias Schnabel, Paul N Bennett, and Marti A Hearst. 2022. SummaC: Re-visiting nli-based models for inconsistency detection in summarization. *Transactions of the Association for Computational Linguistics*, 10:163–177.
- Yixin Liu, Budhaditya Deb, Milagro Teruel, Aaron Halfaker, Dragomir Radev, and Ahmed Hassan. 2023. On improving summarization factual consistency from natural language feedback. In *ACL*.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. In *ACL*.
- Shashi Narayan, Shay Cohen, and Maria Lapata. 2018. Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. In *EMNLP*.
- Artidoro Pagnoni, Vidhisha Balachandran, and Yulia Tsvetkov. 2021. Understanding factuality in abstractive summarization with frank: A benchmark for factuality metrics. In *NAACL*.

- Nicholas Pangakis and Sam Wolken. 2024. Knowledge distillation in automated annotation: Supervised text classification with llm-generated training labels. In *NLP+ CSS*.
- Sanjana Ramprasad, Kundan Krishna, Zachary C Lipton, and Byron C Wallace. 2024. Evaluating the factuality of zero-shot summarizers across varied domains. In *EACL*.
- Mathieu Ravaut, Shafiq Joty, Aixin Sun, and Nancy F Chen. 2023. On context utilization in summarization with large language models. *arXiv e-prints*, pages arXiv–2310.
- Hwanjun Song, Hang Su, Igor Shalyminov, Jason Cai, and Saab Mansour. 2024. FineSurE: Fine-grained summarization evaluation using llms. In *ACL*.
- Liyan Tang, Tanya Goyal, Alexander R Fabbri, Philippe Laban, Jiacheng Xu, Semih Yavuz, Wojciech Kryściński, Justin F Rousseau, and Greg Durrett. 2022. Understanding factual errors in summarization: Errors, summarizers, datasets, error detectors. In *ACL*.
- Liyan Tang, Igor Shalyminov, Amy Wing-mei Wong, Jon Burnsky, Jake W Vincent, Yu’an Yang, Siffi Singh, Song Feng, Hwanjun Song, Hang Su, et al. 2024. TofuEval: Evaluating hallucinations of llms on topic-focused dialogue summarization. In *NAACL*.
- Prasetya Utama, Joshua Bambrick, Nafise Sadat Moosavi, and Iryna Gurevych. 2022. Falsesum: Generating document-level nli examples for recognizing factual inconsistency in summarization. In *ACL*.
- Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020. Asking and answering questions to evaluate the factual consistency of summaries. In *ACL*.
- Shuohang Wang, Yang Liu, Yichong Xu, Chenguang Zhu, and Michael Zeng. 2021. Want to reduce labeling cost? gpt-3 can help. In *EMNLP*.
- Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang, Kathleen Mckeown, and Tatsunori B Hashimoto. 2024. Benchmarking large language models for news summarization. *Transactions of the Association for Computational Linguistics*, 11:39–57.
- Ming Zhong, Da Yin, Tao Yu, Ahmad Zaidi, Mutethia Mutuma, Rahul Jha, Ahmed Hassan, Asli Celikyilmaz, Yang Liu, Xipeng Qiu, et al. 2021. Qmsum: A new benchmark for query-based multi-domain meeting summarization. In *NAACL*.
- Chenguang Zhu, Yang Liu, Jie Mei, and Michael Zeng. 2021. Mediasum: A large-scale media interview dataset for dialogue summarization. In *NAACL*.
- Rongxin Zhu, Jianzhong Qi, and Jey Han Lau. 2023. Annotating and detecting fine-grained factual errors for dialogue summarization. In *ACL*.

Dataset	Source	Annotators	Kappa	Annotation Scheme
AggreFact † (Tang et al., 2022)	CNN/DM XSUM	mixed	-	binary
DiaSumFact (Zhu et al., 2023)	SAMSum QMSum	2 in-house students	0.49	{NoE, EntE, PredE, CirE, CorefE, LinkE, Others}
TofuEval (Tang et al., 2024)	MediaSum MeetingBank	2 expert linguists	0.42 0.34	binary / {extrinsic information error, misreferencing error, stating opinion as fact error, reasoning error, tense/modality error, contradiction error, nuanced meaning shift error, others}
Ramprasad’24 (Ramprasad et al., 2024)	BillSum PubMed	2 expert attorneys 2 expert medical doctors	0.17 0.11	binary / {intrinsic, extrinsic, mixed, others}

Table 6: Summary of the human-labeled datasets. We report Cohen’s kappa in the original by default. In DiaSumFact, we report the average Cohen’s kappa across six annotation groups, each consisting of two annotators. †: we do not report Cohen’s kappa since AggreFact integrates various datasets, some of which include Cohen’s kappa values and others that do not.

	Number of Documents	Number of Summaries	Number of Summarizers	Number of Domains	Doc. Length in Words
Data with Human Feedback	2,499	5,853	17	6	81-2,989 (531)
Data with LLM Feedback (Ours)	10,877	102,640	10	7	5-3,847 (910)

Table 7: Comparison of training datasets with human and LLM feedback. Doc. Length in Words indicates the min-max (average) of the document length in words.

## A Human-labeled Dataset Details

Table 6 summarizes the details of the annotations for each dataset in our human-labeled datasets. The aggregation of these datasets covers various domains and text types. We briefly describe the human-labeled dataset we used.

### A.1 Source Datasets

**AggreFact** (Tang et al., 2022) is a factuality evaluation benchmark that includes two datasets in the news domain; CNN/DM (Hermann et al., 2015) and XSum (Narayan et al., 2018). AggreFact integrates nine datasets from FactCC (Kryściński et al., 2020), Wang’20 (Wang et al., 2020), SummEval (Fabbri et al., 2021), Polytope (Huang et al., 2020), Cao’22 (Cao et al., 2022), XSum-Faith (Maynez et al., 2020), FRANK (Pagnoni et al., 2021), Goyal’21 (Goyal and Durrett, 2021), and CLIFF (Cao and Wang, 2021).

**DiaSumFact** (Zhu et al., 2023) collects fine-grained sentence-level factual error annotations for evaluating dialogue summarization. It spans two dialogue domains: daily conversations, containing chit-chat, and meetings, sourced from SAMSum (Gliwa et al., 2019) and QMSum (Zhong et al., 2021), respectively.

**TofuEval** (Tang et al., 2024) contains two dialogue datasets for benchmarking automated eval-

uators in factuality. It covers two domains, Interview (MediaSum) (Zhu et al., 2021) and Meeting (MeetingBank) (Hu et al., 2023). Each summary is a topic-based summary generated by LLMs and includes sentence-level human annotations for factuality evaluation.

**Ramprasad’24** (Ramprasad et al., 2024) addresses the news domain as well as two specialized domains: medicine (PubMed) (Cohan et al., 2018) and legal (BillSum) (Kornilova and Eidelman, 2019). It releases human annotations from domain experts to assess the factuality of model-generated summaries.

### A.2 Label Consolidation

We aggregate human-labeled data at the sentence level from existing studies, which differ in terms of granularity and annotation schemes when labeling factual consistency. AggreFact and TofuEval provide majority-agreed binary labels for each summary sentence, eliminating the need for consolidation. However, Ramprasad’24 and DiaSumFact consist of two labels for each summary sentence from two annotators, without majority agreement. For these two datasets, if both annotators agreed on ‘no error’, the data was labeled as ‘no error’. If they agreed on ‘error’ it was labeled as ‘error’. If their results differed, the data was not included in our dataset.



Source Dataset	# of Doc.		# of Label 0		# of Label 1		Doc. Length in Words	
	Train	Test	Train	Test	Train	Test	Train	Test
AggreFact	4,130	111	2,325	41	2,754	83	81-2,147 (470.5)	165-1,303 (501.4)
DiaSumFact	339	27	599	43	394	28	93-585 (274.1)	97-342 (238.7)
TofuEval	1,241	531	2,751	1,167	641	322	710-1,199 (963.7)	739-1,185 (919.3)
Ramprasad'24	143	24	421	65	19	6	682-2,989 (1,725.9)	916-2,349 (1,721.9)

Table 8: Statistics of train and test data annotated by humans according to the source. The number of documents (# of Docs) is counted at the document level, while the number of labels (# of Labels) is counted at the sentence level. The 'Doc. Length in Words' column indicates 'min-max (average)' of the document length in words.

Llama-3-8B-Inst. (zero-shot)	Error Type	Pred #	Correct #	Accuracy
		OutE	1,093	101
	EntE	579	168	29.0%
	PredE	101	12	11.9%
	CirE	395	30	7.6%
	GramE	25	6	24.0%
	LinkE	9	0	0.0%
	CorefE	11	2	18.2%

Table 9: Error localization accuracy details in the zero-shot setting. Among the 9 error types, "No Error" and "Other Error" were excluded from the calculation. A total of 2,213 sentences were predicted as errors, of which 319 matched the correct error category.

### A.3 Dataset Split

We aggregate a total of 6,546 document-summary pairs of human feedback data, consisting of 1,772 from TofuEval, 4,241 from AggreFact, 336 from DiaSumFact, and 167 from Ramprasad'24. Then, we split it into 5,853 training set and 693 test set. Data classified as the test set in the original datasets are also included as test data in our dataset. The test set is consistently used as ground truth to evaluate all the models in Table 1. Refer to Table 7 for a comparison between the entire set of human feedback data and LLM feedback data. The breakdown of human feedback data into train and test sets is summarized in Table 8.

### A.4 Testset for Error Localization

We construct an additional test set of 1,286 document-summary pairs from FRANK (Pagnoni et al., 2021), which is a test set tailored for error localization evaluation in Table 2. This dataset consists of labels annotated by three human annotators for each summary sentence across seven error types, identical to those in FineSurE (Song et al., 2024). In sentence-level error localization, if the model-predicted error type matched any one of the three human annotations, it was considered correct. Detailed error localization performance of Table 2 is shown in Tables 9 and 10.

Llama-3-8B-Inst. (fine-tuned)	Error Type	Pred #	Correct #	Accuracy
		OutE	281	80
	EntE	305	160	52.5%
	PredE	42	17	40.5%
	CirE	55	17	30.9%
	GramE	45	10	22.2%
	LinkE	40	8	20.0%
	CorefE	1	0	0.0%

Table 10: Error localization accuracy details in the fine-tuning setting. Among the 9 error types, "No Error" and "Other Error" were excluded from the calculation. A total of 769 sentences were predicted as errors, of which 292 matched the correct error category.

You will receive a document followed by a corresponding summary.  
Your task is to assess the factuality of each summary sentence.

Evaluation Criteria:

**Factual Consistency** - the factual alignment between the summary sentence and the summarized source. A factually consistent summary contains only statements that are entailed by the source document.

Instruction:  
First, compare each summary sentence with the document.  
Second, answer "consistent" or "inconsistent" for each sentence in the summary.

Provide your answer in **JSON format**. The answer should be a list of dictionaries whose keys are "sentence" and "response":  
[{"sentence": "first sentence", "response": "consistent"},  
{"sentence": "second sentence", "response": "inconsistent"},  
{"sentence": "third sentence", "response": "inconsistent"}]

**Document:**  
{input document}

**Summary with N sentences:**  
{summary sentence 1}  
{summary sentence 2}  
...  
{summary sentence N}

**JSON Output:**

Figure 2: Prompt for fact verification ("Binary Label" in Table 3).

## B Fact Verification Prompts

We use three prompts to generate LLM feedback on fact verification, progressively increasing their granularity, as seen in our ablation of Table 3. The first one focuses on fact-checking using binary labels, as shown in Figure 2, while the second adds reasoning to the first one, as shown in Figure 3, and the third further incorporates error localization, as shown in Figure 4.

Specifically, the first prompt asks the LLM to assess the factual consistency of each summary

You will receive a document followed by a corresponding summary. Your task is to assess the factuality of each summary sentence.

Evaluation Criteria:

**Factual Consistency** - the factual alignment between the summary sentence and the summarized source. A factually consistent summary contains only statements that are entailed by the source document.

Instruction:  
 First, compare each summary sentence with the document.  
 Second, provide a single sentence explaining whether the sentence is factually consistent.  
 Third, answer "consistent" or "inconsistent" for each sentence in the summary.

Provide your answer in **JSON format**. The answer should be a list of dictionaries whose keys are "sentence", "reason", and "response":  
 [{"sentence": "first sentence", "reason": "your reason", "response": "consistent"},  
 {"sentence": "second sentence", "reason": "your reason", "response": "inconsistent"},  
 {"sentence": "third sentence", "reason": "your reason", "response": "inconsistent"}]

**Document:**  
 {input document}

**Summary with N sentences:**  
 {summary sentence 1}  
 {summary sentence 2}  
 ...  
 {summary sentence N}

**JSON Output:**

Figure 3: Prompt for fact verification ("Binary Label + Reasoning" in Table 3).

You will receive a document followed by a corresponding summary. Your task is to assess the factuality of each summary sentence across **nine categories**:

- \* no error: the statement aligns explicitly with the content of the document and is factually consistent with it.
- \* out-of-context error: the statement contains information not present in the document.
- \* entity error: the primary arguments (or their attributes) of the predicate are wrong.
- \* predicate error: the predicate in the summary statement is inconsistent with the document.
- \* circumstantial error: the additional information (like location or time) specifying the circumstance around a predicate is wrong.
- \* grammatical error: the grammar of the sentence is so wrong that it becomes meaningless.
- \* coreference error: a pronoun or reference with wrong or non-existing antecedent.
- \* linking error: error in how multiple statements are linked together in the discourse (for example temporal ordering or causal link).
- \* other error: the statement contains any factuality error which is not defined here.

Instruction:  
 First, compare each summary sentence with the document.  
 Second, provide a single sentence explaining which factuality error the sentence has.  
 Third, answer the classified error category for each sentence in the summary.

Provide your answer in **JSON format**. The answer should be a list of dictionaries whose keys are "sentence", "reason", and "category":  
 [{"sentence": "first sentence", "reason": "your reason", "category": "no error"},  
 {"sentence": "second sentence", "reason": "your reason", "category": "out-of-context error"},  
 {"sentence": "third sentence", "reason": "your reason", "category": "entity error"}]

**Document:**  
 {input document}

**Summary with N sentences:**  
 {summary sentence 1}  
 {summary sentence 2}  
 ...  
 {summary sentence N}

**JSON Output:**

Figure 4: Prompt for fact verification ("Binary Label + Reasoning + Error Localization" in Table 3, which is exactly the same with FineSurE (Song et al., 2024)).

sentence against the source document, labeling sentences as either "consistent" or "inconsistent." The second prompt adds complexity by requiring the LLM to not only judge consistency but also provide a brief explanation for each sentence's classification. The third prompt further refines the process by asking the LLM to categorize specific types of factual errors across nine categories, allowing for detailed error identification.

## C Factual Error Types

We follow the error taxonomy suggested by Pagnoni et al. (2021) for feedback generation.

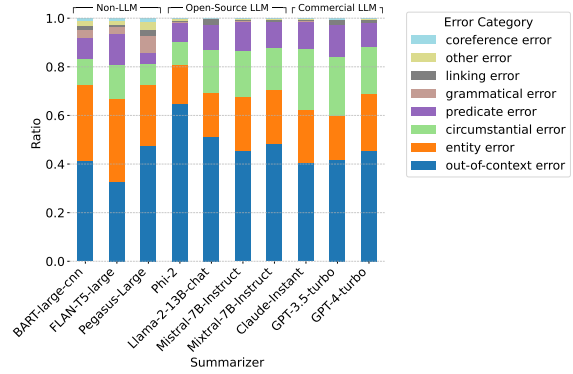


Figure 5: Error category distribution of summaries with LLM feedback for each summarizer, where the error category is estimated using the automated fact verification.

We provide explanations for each error category.

**Out of Context Error (OutE)** indicate that summary statements include information not found in the document, which generally refers to an extrinsic error.

**Entity Error (EntE)** means errors where the core arguments such as subject and object are wrong. This error typically occurs when the generated summary swaps entities.

**Predicate Error (PredE)** refers to errors where the predicate in summary statements is not consistent with the document.

**Circumstance Error (CirE)** occurs when additional information specifying the context around a predicate, such as location, time, or manner, is incorrect.

**Grammatical Error (GramE)** encompasses errors in summary statements where significant grammatical mistakes make them meaningless.

**Discourse Link Error (LinkE)** refers to error in how multiple statements are linked in the discourse, such as incorrect temporal ordering or causal links.

**Coreference Error (CoreE)** is an error where pronouns or references are incorrect antecedents, causing ambiguity.

## D LLM Feedback Data Details

### D.1 Dataset Construction

We generate 102,640 summaries from a total of 10,877 documents, including 18,846 from CNN/DM, 16,744 from MediaSum, 18,990 from DialogSum, 14,427 from MeetingBank, 19,202

Source Dataset	# of Doc.	# of Label 0	# of Label 1	Doc. Length in Words
CNN/DM	18,846	69,215	8,486	54-2,133 (774.8)
MediaSum	16,744	57,334	11,714	96-3,809 (1,546.3)
DialogSum	18,990	36,148	17,545	50-1,084 (187.6)
MeetingBank	14,427	41,187	14,169	126-3,847 (1,164.65)
WikiHow	19,202	31,466	10,587	5-801 (83.7)
GovReport	3,045	14,378	728	180-3,777 (2,775.7)
PubMed	11,386	47,757	3,606	21-3,721 (1,823.5)

Table 11: Statistics of the training data labeled by LLM (Llama-3-70B-Inst.) according to the source. The number of documents (# of Docs) is counted at the document level, while that of labels (# of Labels) is counted at the sentence level. The 'Doc. Length in Words' column indicates 'min-max (average)' of the document length in words.

Model Name	HuggingFace/API Checkpoints
<b>Summary Generation</b>	
BART <sub>large-cnn</sub>	facebook/bart-large-cnn
FLAN-T5 <sub>large</sub>	google/flan-t5-large
Pegasus <sub>large</sub>	google/pegasus-large
Phi-2	microsoft/phi-2
Llama-2 <sub>13B-chat</sub>	meta-llama/Llama-2-13b-chat-hf
Mistral <sub>7B-Instruct</sub>	mistralai/Mistral-7B-Instruct-v0.1
Mixtral <sub>8x7B-Instruct</sub>	mistralai/Mixtral-8x7B-Instruct-v0.1
Claude <sub>Instant</sub>	claude-instant-1.2
GPT-3.5 <sub>turbo</sub>	gpt-3.5-turbo-0125
GPT-4 <sub>turbo</sub>	gpt-4-turbo-2024-04-09
<b>Feedback Generation</b>	
Llama-3 <sub>70B-Instruct</sub>	meta-llama/Meta-Llama-3-70B-Instruct
<b>Fine-tuning</b>	
Llama-3 <sub>8B-Instruct</sub>	meta-llama/Meta-Llama-3-8B-Instruct

Table 12: The model checkpoints.

from WikiHow, 3,045 from GovReport, and 11,386 from PubMed. The summaries and feedback are generated by 10 different language models and Llama-3-70B-Instruct, respectively. This dataset is exclusively used as a training set for fine-tuning the model, not as a test set. We provide the statistics of the LLM feedback dataset in Table 11. Table 12 provides the details of the experiment’s models.

## D.2 Error Category Distribution

We analyze the distribution of error categories based on feedback provided by Llama-3-70B-Instruct, evaluating summaries generated by 10 different LLMs. As shown in Table 13, the number of errors decreases as we move from non-LLMs to open-source LLMs and then to commercial LLMs. Additionally, we find that most summarizers exhibit a higher proportion of out-of-context errors and entity errors, while coreference errors are the least frequent. We provide the summary error type distribution in Figure 5.

## E Agreement with Humans per Domain

As shown in Table 14, the performance across the News, Interview, and Meeting domains reveals varying levels of agreement with human judgments in fact verification, with each domain

Summarizer	No Error	Error	Error Ratio
<b>Non-LLM</b>			
BART <sub>large-cnn</sub>	25,629	10,310	22.29%
FLAN-T5 <sub>large</sub>	16,360	5,738	20.61%
Pegasus <sub>large</sub>	18,850	3,737	14.20%
<b>Open-Source LLM</b>			
Phi-2	10,917	11,559	33.96%
Llama-2 <sub>13B-chat</sub>	20,040	4,618	15.77%
Mistral <sub>7B-Instruct</sub>	36,979	7,780	14.81%
Mixtral <sub>7B-Instruct</sub>	36,608	9,717	17.34%
<b>Commercial LLM</b>			
Claude <sub>Instant</sub>	42,054	4,771	9.25%
GPT-3.5 <sub>turbo</sub>	41,405	3,631	7.46%
GPT-4 <sub>turbo</sub>	40,791	2,929	6.28%

Table 13: Error Ratio according to summarizers, indicates the proportion of summaries generated by each summarizer that are identified as errors by the feedback generator (Llama-3-70B-Instruct).

presenting unique challenges and insights.

**News** Both QA-based and LLM-based methods showed high agreement with human judgments, demonstrating their effectiveness in handling structured, fact-dense content typically found in news articles.

**Interview** LLM Feedback performed notably well, while QA-based and NLI-based methods struggled, underscoring the difficulties posed by the unstructured and conversational format of interview content.

**Meeting** The results were similar to those in the Interview domain, with the LLM Feedback method outperforming others. However, the overall performances of each type were lower, reflecting the inherent complexity in summarizing and verifying content from meetings.

LLM-based methods stood out for their consistent performance across different domains. This robustness can be attributed to their fine-tuning with aggregated datasets that span a wide variety

Type	Method	News		Interview		Meeting	
		Summ.	Sys.	Summ.	Sys.	Summ.	Sys.
QA-based	QAFactEval	0.614*	0.886*	0.406*	-0.257	0.382*	-0.167
NLI-based	SummaC-Conv	0.515*	0.683*	0.208*	-0.086	0.168*	-0.433
LLM-based	Zero-shot	0.402*	0.605*	0.198*	0.829*	0.181*	0.150
	Human Feedback	0.287*	0.560*	0.443*	0.600	0.468*	0.083
	LLM Feedback	0.573*	0.832*	0.528*	0.886*	0.529*	0.633

Table 14: Agreement with human judgments in fact verification on test data across three domains: News, Interview, and Meeting. The agreement was measured at the summary level using Pearson correlation (Summ.) and at the system level using rank correlation (Sys.) with a significance threshold of p-value < 0.05 (\*). Results are reported per domain only when the test examples exceed 20. Domains with insufficient data, specifically **Daily**, **Legal**, and **Medicine** were excluded due to inflated p-values and statistically insignificant results.

of domains, enabling them to generalize effectively across different types of content.

## F Metrics

We follow the same settings as those presented in recent studies (Song et al. 2024, Liu et al. 2023) to assess the model’s performance and ensure alignment with human judgment.

**bAcc** *Balanced accuracy (bAcc)* is used to address class imbalance when summarizing the performance of a model in a classification task. During sentence-level evaluation, human annotations and fine-tuned LLM classify factual correctness as '0' (No error) and incorrectness as '1' (Error). The formula for bAcc is as follows:

$$\text{bAcc} = \frac{\text{TPR} + \text{TNR}}{2} \quad (1)$$

TPR (True Positive Rate), measures the proportion of correct positive predictions made by the fine-tuned LLM. TNR (True Negative Rate) quantifies the proportion of correct negative predictions made.

**Faithfulness score** For the summary-level and system-level evaluations, the percentage score of faithfulness enables us to assess summaries by aggregating sentence-level fact checks. Let us assume  $S_i = \{s_{i,1}, \dots, s_{i,N}\}$  represents the  $i$ -th summary passage, consisting of  $N$  sentences, where  $s_{i,j}$  denotes the  $j$ -th sentence in the  $i$ -th summary passage. Additionally, let  $S_{i,\text{fact}} \subseteq S_i$  represent the subset of sentences identified as factually correct within this summary. The percentage score of faithfulness for  $S_i$ , with respect to the original document  $D_i$ , is computed as follows:

$$F(D_i, S_i) = \frac{|S_{i,\text{fact}}|}{|S_i|} \quad (2)$$

**Summary-level correlation** To compute the *summary-level correlation*, we define  $F_{gt}$  and  $F_{pred}$  as the faithfulness scores of the ground truth and the prediction, respectively. Let  $D = \{D_1, \dots, D_k\}$  represent the set of input documents, and  $S = \{S_1, \dots, S_k\}$  denote the corresponding set of summaries for these documents. Then the summary-level correlation is computed as follows:

$$\text{Pearson}([F_{gt}(D_1, S_1), \dots, F_{gt}(D_k, S_k)], [F_{pred}(D_1, S_1), \dots, F_{pred}(D_k, S_k)]) \quad (3)$$

**System-level rank correlation** To compute the *system-level rank correlation*, we define  $\mathbf{F}_m = \{F_m(D_1, S_1), \dots, F_m(D_k, S_k)\}$  as the set of percentage scores obtained from the labels given by the summarization model  $m$ . Then, we construct a list of the average percentage scores for all summarization models, denoted as  $[\bar{\mathbf{F}}_{m_1}, \bar{\mathbf{F}}_{m_2}, \dots]$  where,  $\bar{\mathbf{F}}_{m_i} = \frac{1}{|m_i|} \sum_{j=1}^{|m_i|} F_{m_i}(D_j, S_j)$ . Using this list and the Rank function, we create the list  $[\text{rank}_{m_1}, \text{rank}_{m_2}, \dots]$ , where  $\text{rank}_m$  represents the rank of model  $m$ . By the same mechanism, we construct the ground truth list of ranks  $[\text{rank}_{m_1}^*, \text{rank}_{m_2}^*, \dots]$  using the human labels. Finally, the summary-level correlation is computed as follows:

$$\text{Spearman}([\text{rank}_{m_1}, \text{rank}_{m_2}, \dots], [\text{rank}_{m_1}^*, \text{rank}_{m_2}^*, \dots]) \quad (4)$$

The summary-level correlation indicates the agreement between human judgments and LLM, while the system-level rank correlation measures how closely the model rankings align with those provided by humans across various summarizers.

	<p>You will receive a document followed by a corresponding summary.          Your task is to assess the factuality of each summary sentence across <b>nine categories</b>:</p> <ul style="list-style-type: none"> <li>* no error: the statement aligns explicitly with the content of the document and is factually consistent with it.</li> <li>* out-of-context error: the statement contains information not present in the document.</li> <li>* entity error: the primary arguments (or their attributes) of the predicate are wrong.</li> <li>* predicate error: the predicate in the summary statement is inconsistent with the document.</li> <li>* circumstantial error: the additional information (like location or time) specifying the circumstance around a predicate is wrong.</li> <li>* grammatical error: the grammar of the sentence is so wrong that it becomes meaningless.</li> <li>* coreference error: a pronoun or reference with wrong or non-existing antecedent.</li> <li>* linking error: error in how multiple statements are linked together in the discourse (for example temporal ordering or causal link).</li> <li>* other error: the statement contains any factuality error which is not defined here.</li> </ul> <p><b>Instruction:</b>          First, compare each summary sentence with the document.          Second, provide a single sentence explaining which factuality error the sentence has.          Third, answer the classified error category for each sentence in the summary.</p> <p>Provide your answer in <b>JSON format</b>. The answer should be a list of dictionaries whose keys are "sentence", "reason", and "category":</p> <pre>[{"sentence": "first sentence", "reason": "your reason", "category": "no error"}, {"sentence": "second sentence", "reason": "your reason", "category": "out-of-context error"}, {"sentence": "third sentence", "reason": "your reason", "category": "entity error"}]</pre> <p><b>Input</b></p> <p><b>Document:</b>          Chronic condition management has been structured around models that aim at facilitating and development of care through initiatives of self - management , organisational restructuring and information system enhancement . uk s health care framework for chronic condition care draws upon the chronic care model and innovation care for chronic conditions framework . however , consideration to the clarity or vagueness of boundaries between different hierarchical structures and awareness of each structure about its focus and contribution to overall health system are not taken into account by these models . the proposed chronic condition management model is different as it is targeted at primary care which plays a salient role of gate keeping . change implemented here would lead to better clinical effectiveness , increase in efficiency of the system and reduction in cost . this model does not directly call for policy change ; instead it recommends incorporating available evidence into the policy . it calls for empowering the patient by providing health literacy and improving awareness , and for collaboration and cooperation between primary , secondary , tertiary , community and social care services . initiatives towards chronic condition management should be enacted at primary care due to its gate keeping role in the uk s health care system .</p> <p><b>Summary with 3 sentences:</b>          [1] The proposed chronic condition management model emphasizes the role of primary care in the UK’s health care system for managing chronic conditions.          [2] This model advises incorporating available evidence into policies, empowering patients with health literacy and awareness, and fostering collaboration between primary, secondary, tertiary, community, and social care services.          [3] The model does not call for direct policy changes, but rather for addressing the clarity or vagueness of boundaries between hierarchical structures and making each structure aware of its focus and contribution to the overall health system.</p> <p><b>JSON Output:</b></p>
<p><b>Output</b></p>	<pre>[ {"sentence": "The proposed chronic condition management model emphasizes the role of primary care in the UK’s health care system for managing chronic conditions.", "reason": "The sentence accurately summarizes the document’s information about the proposed model’s focus on primary care.", "category": "no error"}, {"sentence": "This model advises incorporating available evidence into policies, empowering patients with health literacy and awareness, and fostering collaboration between primary, secondary, tertiary, community, and social care services.", "reason": "The sentence accurately summarizes the document’s information about the model’s recommendations.", "category": "no error"}, {"sentence": "The model does not call for direct policy changes, but rather for addressing the clarity or vagueness of boundaries between hierarchical structures and making each structure aware of its focus and contribution to the overall health system.", "reason": "The sentence incorrectly states that the model addresses the clarity or vagueness of boundaries between hierarchical structures, which is a criticism of existing models, not a feature of the proposed model.", "category": "out-of-context error"} ]</pre>

Table 15: An example of an input (user prompt) and its corresponding output (assistant prompt) used to fine-tune the Llama-3-8B-Instruct model.