

David vs. Goliath: Can small models leverage LLMs for summarization?

Anonymous ACL submission

Abstract

Recent studies indicate a preference for summaries generated using large language models (LLMs) over those using classical models, highlighting a performance discrepancy. This study explores strategies to narrow the gap between the summaries generated through these two models. To address this, we introduce a novel framework that uses LLM-generated summaries to train classical models, adopting a two-stage training approach to enhance their summary quality. Although classical models are relatively smaller in size than LLMs, through automatic metrics and human evaluations, we can demonstrate that the performances of classical models, trained using LLM-generated references can catch up with LLMs. Our findings create a simple yet potential way to improve classical summarization models by leveraging LLMs. Additionally, we contribute a new dataset **GXSum**¹, enabling further research and promoting development progress in this subject.

1 Introduction

Text summarization plays a pivotal role in the field of natural language processing by condensing articles into concise versions that capture the main information. With the rapid development of deep learning, automatic text summarization systems have made significant progress. (Nallapati et al., 2016a; Vaswani et al., 2017; Li et al., 2018; Shi et al., 2021). More recently, large language models (LLMs) have revolutionized the field of natural language processing. These models exhibit remarkable results in summarization accuracy, particularly under zero-shot and few-shot fine-tuning scenarios (Wang et al., 2023; Basyal and Sanghvi, 2023; Ahmed and Devanbu, 2023). Unlike classical models, LLMs leverage reinforcement learning from human feedback (RLHF) (Kirk et al., 2023),

fine-tuning their outputs to align more closely with human preferences, thereby widening the performance gap with classical models (Wang et al., 2023; Zhang et al., 2024; Fabbri et al., 2021). Some studies even indicate that humans might prefer LLM-generated summaries to those written (or selected) by humans (Liu et al., 2023b,a).

Sweeping over previous research on text summarization, most studies mainly concentrated on developing novel model architectures (Dou et al., 2021; Wang et al., 2022a; Liu et al., 2022) or training method (Stiennon et al., 2020). These efforts improve performance on specific benchmarks, yet they often increase model complexity or compromise training efficiency. However, these efforts still do not bridge the performance gap with LLMs.

Knowledge distillation is a simple and straightforward way to transfer model capabilities from one model to another. To move beyond LLMs in a simple and cost-effective way, we present a two-stage training framework that is expected to allow classical summarization models to rival the performance of LLMs based on the fundamental philosophy of knowledge distillation in this study. More specifically, in the first stage, we leverage LLMs to generate summaries and form a new dataset. Next, we train classical models referring to the new ground truths with the traditional maximum likelihood objective. By doing so, the classical model is expected to not only inherit the advantages of LLMs but also retain the abilities of original designs, delivering better results than LLMs.

In sum, our key contributions are at least three-fold. First, we propose a simple yet efficient framework to enhance the performance of classical models and catch up with LLMs. Second, a series of experiments were used to show that significant performance gains are achievable even with limited data for fine-tuning. Of course, as always, more data yields better results. Third, a new dataset **GXSum** is released to facilitate further research,

*Equal contribution.

¹<https://github.com/anonymous>

| | | | |
|-----|--|---|-----|
| 081 | perform fair comparisons, make results producible, | of LLMs in summarization tasks. | 131 |
| 082 | and promote research progress in the line of re- | | |
| 083 | search. | | |
| 084 | 2 Related Work | 3 LLM-Guided Summarization | 132 |
| 085 | Previous research has demonstrated the exceptional | 3.1 Models | 133 |
| 086 | proficiency of LLMs in generating summaries, out- | In this study, we selected the most advanced Chat- | 134 |
| 087 | performing classical models in both automated | GPT ² provided by OpenAI as an example. To | 135 |
| 088 | evaluation metrics and human assessments. Ad- | minimize the randomness of generated results, we | 136 |
| 089 | ditionally, summaries generated using LLMs, es- | set the temperature parameter of the model to 0, | 137 |
| 090 | pecially in the news domain, have been shown to | whereas other parameters are at their default val- | 138 |
| 091 | be at par with, or even superior to, those crafted | ues to ensure stability and reproducibility of the | 139 |
| 092 | by humans. These results reveal significant potential | experimental results. | 140 |
| 093 | for LLMs on the text summarization task (Victor | For a comprehensive analysis, BART (Lewis | 141 |
| 094 | et al., 2022; Wang et al., 2022b; Goyal et al., 2022). | et al., 2020), PEGASUS (Zhang et al., 2020), and | 142 |
| 095 | Some studies further emphasize that the field of | BRIO (Liu et al., 2022) were chosen as the basic | 143 |
| 096 | summarization is undergoing significant changes, | classic summarization models for our exper- | 144 |
| 097 | suggesting a pivotal moment in summarization re- | iments. These models have been proven in pre- | 145 |
| 098 | search. A thought-provoking question is whether | vious research to possess excellent text summa- | 146 |
| 099 | those human-generated ground truths bound the | rization capabilities, each representing various re- | 147 |
| 100 | performances of classical summarization models | search directions in the field of summarization. The | 148 |
| 101 | (Pu et al., 2023; Zhang et al., 2024). | pre-trained models of BART and PEGASUS are | 149 |
| 102 | The feasibility of using LLMs for generating | sourced from the Transformers Library (Wolf et al., | 150 |
| 103 | source data has been extensively explored. Some re- | 2020), whereas the weights for BRIO are obtained | 151 |
| 104 | search has introduced methods for distilling LLMs | from the GitHub repository of the original paper. | 152 |
| 105 | and employing them in data augmentation tasks | 3.2 Human Referenced Datasets | 153 |
| 106 | (Wang et al., 2021; Ding et al., 2023; Kang et al., | In this study, we adopted two key news summa- | 154 |
| 107 | 2023). Specifically, these methods focus on extract- | rization datasets that are widely used in the re- | 155 |
| 108 | ing the most relevant information from LLMs to | search of summarization models and the evalua- | 156 |
| 109 | enrich training datasets, thereby enhancing model | tion of large language model performance: the | 157 |
| 110 | performance without the need for extensive com- | Extreme Summarization Dataset (abbreviated as | 158 |
| 111 | putational resources. Notably, a series of studies | XSum) (Narayan et al., 2018) ³ and the CNN / | 159 |
| 112 | have demonstrated the use of LLMs to generate | DailyMail News Summarization Dataset (abbrev- | 160 |
| 113 | both final answers and task-related descriptions, | iated as CNNDM) ⁴ (Nallapati et al., 2016b). The | 161 |
| 114 | which aid in training smaller models for reasoning | XSum dataset is comprised of press releases from | 162 |
| 115 | tasks (Li et al., 2022; Shridhar et al., 2023; Hsieh | the British Broadcasting Corporation, whereas the | 163 |
| 116 | et al., 2023). In the realm of text summarization, | CNNDM dataset compiles news articles from the | 164 |
| 117 | Wang et al. (2021) have used GPT-3 (Brown et al., | Cable News Network (CNN) and the Daily Mail. | 165 |
| 118 | 2020) to generate reference summaries. Concur- | Notably, these two datasets differ significantly in | 166 |
| 119 | rently, Gekhman et al. (2023) proposed the use of | their nature. Compared to CNNDM, the summary | 167 |
| 120 | LLMs for annotating summary factual consistency | reference texts in XSum mostly contain only one | 168 |
| 121 | (Maynez et al., 2020), facilitating the training of | to two sentences, posing a significant challenge for | 169 |
| 122 | models to evaluate factual consistency. Moreover, | summarization models to refine and extract core | 170 |
| 123 | Liu et al. (2023c) have explored further fine-tuning | information for the summary. Table 1 shows the | 171 |
| 124 | of news summaries generated by the GPT series | ROUGE scores (cf. section 4.2) of classic models | 172 |
| 125 | for the summarization domain. | on the XSum and CNNDM datasets. | 173 |
| 126 | Therefore, in this paper, we expand the dataset | | |
| 127 | and thoroughly analyze the differences between | | |
| 128 | LLMs and human summarization. In the subse- | | |
| 129 | quent research, we will further train the summaries | | |
| 130 | generated using LLMs, aiming to redefine the role | | |

²GPT-4-Turbo (gpt-4-1106-review)

<https://platform.openai.com/docs/models/overview>

³<https://github.com/EdinburghNLP/XSum>

⁴<https://cs.nyu.edu/~kcho/DMQA/>

| Models | XSum | | | CNNDM | | |
|---------|-------|-------|-------|-------|-------|-------|
| | R-1 | R-2 | R-L | R-1 | R-2 | R-L |
| BART | 45.14 | 22.27 | 37.25 | 44.16 | 21.28 | 40.90 |
| PEGASUS | 47.21 | 24.56 | 39.25 | 44.17 | 21.47 | 41.11 |
| BRIO | 49.07 | 25.59 | 40.40 | 47.78 | 23.55 | 44.57 |

Table 1: BART, PEGASUS, and BRIO’s ROUGE scores on the XSum and CNNDM datasets.

3.3 LLM Referenced Dataset

As one of the core objectives of our research, we created a dataset comprising summaries generated by LLMs to serve as reference summaries. This dataset is based on XSum and CNNDM, maintaining the format of the original datasets. To leverage the ChatGPT API for generating high-quality summaries, we have meticulously designed a prompt template that specifically emphasizes the role of ChatGPT as a summary writer. Additionally, to better control the summary length, we included a description of the length limit as a soft constraint in the prompt and set the API `max_tokens` parameter as a hard constraint. The detailed design of the prompt is presented in Appendix A.1. For the source text, we designated the document from XSum and the article from CNNDM as the variables. During the summary generation process, the length restriction was set to ensure that the difference in lengths between the newly generated summaries and the original reference summaries remained within a range of plus or minus five tokens. We provide an example of our summary generation process in Appendix A.2.

3.4 Implementation Details

3.4.1 Data Processing

We extracted a sample comprising 20,000 data points from the training set and 1,100 data points from the validation set. These samples were subjected to the LLM summarization workflow to produce reference summaries. This subset was designated as the *Small* variant. In contrast, the test set underwent comprehensive processing to guarantee a robust and reliable evaluation. Data processing was conducted on both the XSum and CNNDM datasets to ensure uniformity and accuracy in our analyses.

3.4.2 Training Details

The initiation of training for each model leveraged checkpoints that had been previously fine-tuned on the benchmarked XSum and CNNDM datasets. These fine-tuned checkpoints used for BART⁵, PEGASUS⁶ and BRIO⁷ were obtained from the Huggingface library. For optimization, the AdamW optimizer was employed, incorporating a weight decay of 0.01 and an initial learning rate of 0.00002. A linear learning rate scheduler was applied without any warm-up steps. Model performance evaluation on the validation set informed the selection of checkpoints, whereas performance metrics on the test set were documented and reported.

3.5 Evaluation Methods

To validate the performance of our model, we use two primary evaluation methods: human validation and automatic metrics. Initially, human validation gauges the summaries’ quality from readers’ viewpoints. Automatic metrics are used to determine whether the fine-tuning process is functioning properly and toward the training objectives.

3.5.1 Human Evaluation Protocol

As the main evaluation methods of this study, we adopted three common forms of human validation, including the Likert scale scoring, pairwise comparison, and multiple candidate ranking.

The Likert scale scoring is the most used method in human validation assessments. The evaluation process involves presenting a source text and its corresponding generated summary, where human annotators are required to score the summary on several aspects of performance. In this research, we defined five distinct aspects for evaluation: relevance, consistency, fluency, coherence, and infor-

⁵<https://huggingface.co/facebook/bart-large-xsum> and <https://huggingface.co/facebook/bart-large-cnn>

⁶<https://huggingface.co/google/pegasus-xsum> and https://huggingface.co/google/pegasus-cnn_dailymail

⁷<https://huggingface.co/Yale-LILY/brio-xsum-cased> and <https://huggingface.co/Yale-LILY/brio-cnndm-cased>

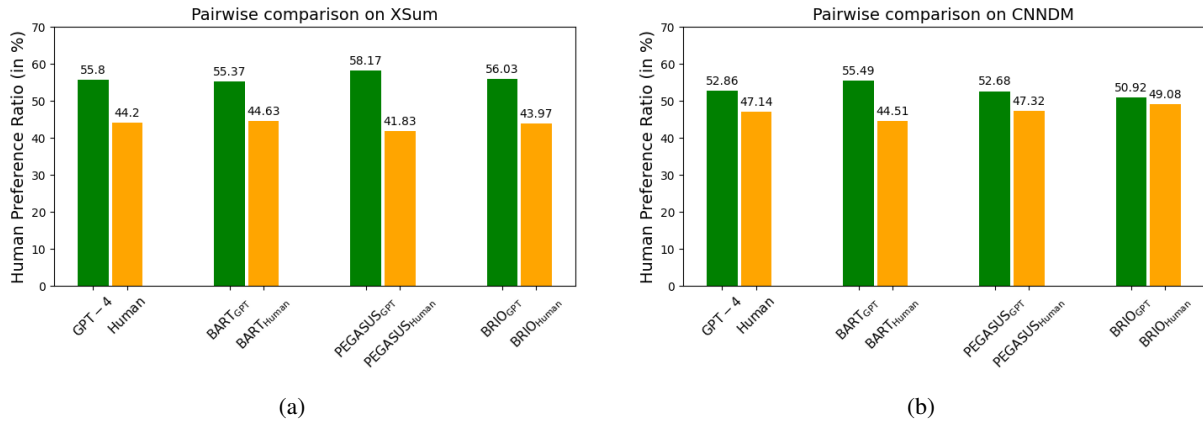


Figure 1: Pairwise Comparison on XSum and CNNDM

mativeness. Detailed guidelines for these metrics are elaborated in Appendix B.1. Through these metrics, human annotators can more comprehensively score the overall quality of summaries. The scoring range is set from 1 (worst) to 5 (best).

Pairwise comparison is a human validation evaluation method based on relative comparison. Given a source text and two summaries generated by different models, assessors are asked to select the one with the better quality.

Multiple candidate rating is an advanced and complex variation of the pairwise comparison method. Assessors are compelled to examine a set of summaries for a given source text and assign a unique rating to each, reflecting the overall quality of each summary. Therefore, the method facilitates a thorough evaluation of the performance variations across various summarization models. Within our experiment, we established a rating scale from 1 (lowest quality) to 5 (highest quality).

3.5.2 Automatic Evaluation Metrics

We adopted Recall-Oriented Understudy for Gisting Evaluation (ROUGE) (Lin, 2004) as our automatic evaluation metric for summarization effectiveness. ROUGE is crucial in performing summarization research, serving as a standard for comparing the similarity and quality between computer-generated and human-crafted reference summaries. This study employs three ROUGE variants: ROUGE-1 (R-1), ROUGE-2 (R-2), and ROUGE-L (R-L). ROUGE-1 assesses unigram similarity to gauge informational content. ROUGE-2 evaluates bigram similarity for fluency. ROUGE-L focuses on the longest common subsequence to determine core content extraction.

4 Experiment Result and Analysis

4.1 Human preference

The collection of human annotations contains evaluations of summaries generated by models that were fine-tuned on the *Small* dataset. These evaluations were obtained through a combination of crowd-sourced contributors and expert judgments. **Crowd-Sourced Annotations** The annotation data from crowd-sourced participants were meticulously gathered via the Amazon Mechanical Turk (MTurk) platform, the detailed recruitment setting is described in Appendix B.2. We collected annotations for a sample of 1,000 articles from the XSum and CNNDM test sets, which were summarized using various pairs of systems. To establish a baseline for comparison, we used models fine-tuned on the original datasets. These baseline models were then compared to those fine-tuned on the LLM-reference dataset. Each summary evaluation will be conducted by three separate crowd annotators, using the Likert scale scoring and pairwise comparison methods detailed in Section 3.5.1.

Figures 1b and 1a depict the crowd-sourced winning rates from the pairwise comparisons for each dataset. Systems trained using human references are denoted with the subscript *Human*, whereas those trained using GPT-4 references are marked with the subscript *GPT*. Pairwise comparisons were conducted between these two training settings, employing the same base model for each comparison. We made the following observations:

(1) The summaries used as training references, generated by the GPT-4 system, surpassed those written by humans in terms of human preference on both the XSum and CNNDM tasks. This outcome substantially supports to the hypotheses posited

| Dataset | System | Relevance | Consistency | Fluency | Coherence | Informativeness |
|---------|------------------------|-----------|-------------|---------|-----------|-----------------|
| | Human Base | 0% | 0% | 0% | 0% | 0% |
| XSum | GPT-4 | +13.8% | +13.2% | +9.3% | +7.5% | +3.6% |
| | BART _{GPT} | +17% | +15.5% | +10.9% | +11.3% | +4.2% |
| | PEGASUS _{GPT} | +18.3% | +15.4% | +14.5% | +16.5% | +7.4% |
| | BRIO _{GPT} | +11% | +8.3% | +9% | +7% | +3.3% |
| CNNDM | GPT-4 | +3.58% | +1.6% | +5.6% | +1.2% | -0.2% |
| | BART _{GPT} | +0.2% | +0.7% | +1.4% | +1.4% | +0.9% |
| | PEGASUS _{GPT} | -1.1% | +3.1% | +1.5% | +1.8% | +1.4% |
| | BRIO _{GPT} | -1.9% | +2.9% | +0.7% | +0.9% | -0.5% |

Table 2: Evaluation through Crowd-Sourced Likert Scale Scoring, which models referenced by humans serve as the baseline for comparison (default as 0%). The report highlights the percentage difference in occurrences where one system is adjudged to *outperform* the other. For instance, GPT-4 exceeds human writers in Relevance by 13.8% on the XSum dataset. In case of a tie, both systems are recognized as winners.

in the related work (Goyal et al., 2022; Liu et al., 2023c; Pu et al., 2023).

(2) **The model trained using the references generated by GPT-4 consistently outperformed the model trained using human-generated references**, demonstrating that supervised training with high-quality references can enhance summarization model performance. On examining the discrepancies between datasets, it was observed that, in comparison to XSum, the performance advantage of the system guided by GPT-4 on CNNDM was less marked.

The observed performance discrepancies, as revealed through pairwise comparisons, prompted us to conduct an in-depth analysis across various systems and datasets, using the Likert scale scoring to quantify these differences with better accuracy. Table 2 illustrates the comparative performance of models referenced by humans versus those referenced by GPT-4, using the Likert scale scoring. From these results, we observe that:

(1) In the XSum dataset, GPT-4 referenced models outperformed across all metrics. Detailed analysis showed the most significant improvement in summary relevance with GPT-4 training, although informativeness remained the weakest point as per the annotators. This is attributed to the XSum requirement (Narayan et al., 2018) for highly abstract, single-sentence summaries, leading to inevitable information loss in both human-generated and GPT-4-generated summaries, with the human ones more susceptible to bias or misinformation. Our experimental results corroborate that GPT-4 can effectively improve these issues. For a more detailed

case study, refer to the Appendix C.

(2) Our analysis of the CNNDM dataset reveals that although GPT-4 referenced summaries still outshine human-generated ones on several metrics, the margin of advantage has significantly narrowed. The Performance of the models trained on GPT-4 and human references are closely matched, with discrepancies often within a 1-2% margin. This trend is likely influenced by the CNNDM dataset approach (Nallapati et al., 2016b) of collecting human reference summaries, which involves compiling human-written summary bullets from CNN and Daily Mail articles. These summaries lean toward an extractive nature, closely mirroring the original article content, and their length helps minimize information omission. Thus, the inherent advantages of GPT-4 training are less pronounced in this context.

The results derived from crowd-sourced annotations significantly validate our approach of using LLM-guided training sets. However, potential reliability concerns of our experiment exist because of the variability in nonexpert summary judgments, as highlighted by prior studies (Callison-Burch and Dredze, 2010; Goyal et al., 2022; Zhang et al., 2024). This variability, reflecting the subjective nature of human evaluation and minor performance disparities between systems, can impact the consistency of annotation quality. To address this, we conducted additional analyses with expert reviewers to ensure more dependable evaluations.

Expert Annotations To ensure the rigor of expert analysis, we established specific criteria for the selection of annotators, focusing on those with a

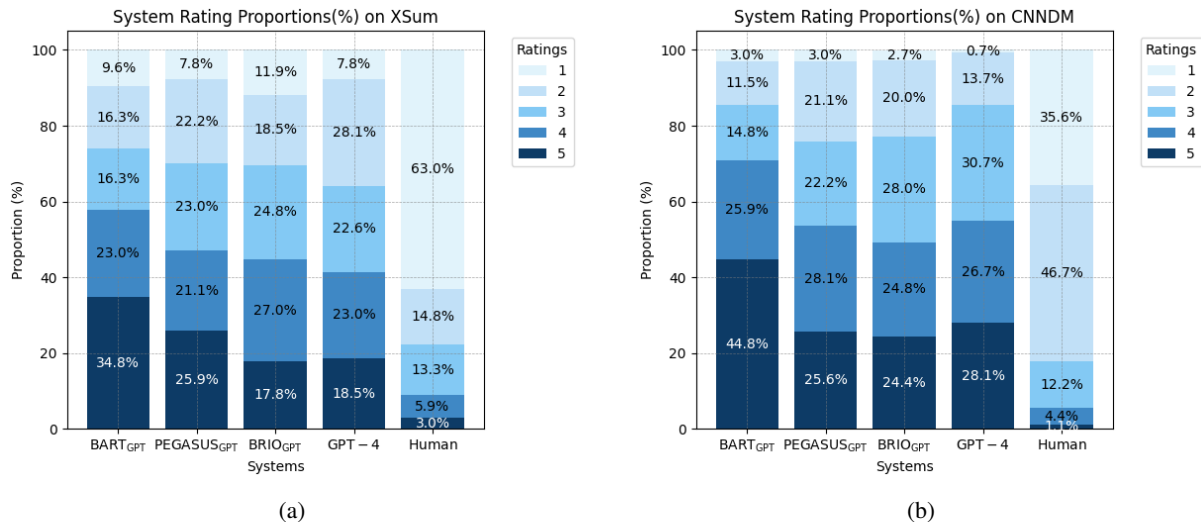


Figure 2: Rating proportions in XSum and CNNDM

requisite level of expertise. We collected annotations for a sample of 100 articles from the XSum and CNNDM test sets. The evaluation of each summary was entrusted to three distinct expert annotators who applied the Multiple Candidates Rating Methods as delineated in Section 3.5.1. Additionally, annotators were required to provide reviews of their annotations, enabling verification of results. The candidates the position of an expert annotator are hired from the Upwork platform. The detailed recruitment setting is described in Appendix B.3.

Figures 2b and 2a illustrate the rating distributions (1-5) for each system according to expert evaluations. The analysis yields two key insights: (1) Expert raters show a clear preference for summaries generated using GPT-4 and GPT-4-assisted systems over those written by humans. This supports our hypothesis based on crowd-sourced annotations, confirming the ability of our system to produce summaries aligned more with human preferences.

(2) Significantly, **models trained using the GPT-4 references achieve, and sometimes surpass, the performance of GPT-4 in expert assessments**, achieving a 68% inter-annotators agreement. This indicates that using our training methodology, smaller models can reach the efficacy of LLMs.

To enhance understanding of our findings, we delved into detailed case studies, referenced in Appendix C. The analysis of annotators' reviews revealed a preference for our fine-tuned models, which produce summaries containing significantly more relevant information and demonstrating better

coherence.

4.2 Automatic Metric

Next, in Table 3, we compare various summary generation models on the XSum and CNNDM datasets, using ROUGE scores for evaluation. This analysis contrasts human-generated summaries with those generated from GPT-4, noting lower ROUGE scores when comparing GPT-4 outputs to human references, highlighting differences in sentence structure and expression. Our results indicate variability in model performance, with GPT's BRIO model leading in ROUGE-1 and ROUGE-L scores on CNNDM, and GPT-based models surpassing human performance on XSum in these scores. Despite this, a significant performance gap exists between the best models and human summaries, particularly on XSum's ROUGE-2 scores. This result shows the strength of GPT-based models in abstract text generation, despite the challenges in closely mimicking human summarization.

5 Comparative Study

To facilitate a more comprehensive analysis of our training methodology, we present a comparative experiment from various perspectives in this section.

5.1 Training Efficiency

In Section 4.2, we detail the ROUGE score performance of various systems fine-tuned on a dataset of 20,000 GPT-4 generated references. The results show a discernible performance gap between the ROUGE scores achieved by our model and those reported in the original papers (Lewis et al., 2020;

| Reference | Hypothesis | XSum | | | CNNDM | | |
|--------------|------------------------|-------|-------|-------|-------|-------|-------|
| | | R-1 | R-2 | R-L | R-1 | R-2 | R-L |
| GPT-4 | Human | 24.95 | 5.64 | 18.59 | 36.80 | 10.89 | 31.91 |
| | BART _{GPT} | 45.36 | 19.59 | 36.28 | 48.97 | 20.84 | 41.03 |
| | PEGASUS _{GPT} | 43.71 | 18.68 | 35.07 | 46.28 | 20.54 | 39.10 |
| | BRIO _{GPT} | 47.37 | 21.30 | 38.55 | 50.03 | 21.96 | 41.73 |
| Human | GPT-4 | 24.95 | 5.64 | 18.57 | 36.80 | 10.90 | 32.05 |
| | BART _{GPT} | 26.39 | 6.61 | 19.10 | 40.05 | 14.86 | 35.08 |
| | PEGASUS _{GPT} | 28.00 | 7.94 | 20.77 | 40.50 | 16.18 | 35.76 |
| | BRIO _{GPT} | 26.81 | 7.01 | 19.81 | 40.39 | 15.19 | 35.20 |

Table 3: Evaluation of ROUGE Scores after Fine-Tuning with 20,000 GPT-4 Summaries. This table presents the calculated ROUGE scores, comparing various **Hypotheses** with **References**.

Zhang et al., 2020; Liu et al., 2022), particularly concerning the XSum dataset. Therefore, we questioned whether fine-tuning the model on a larger dataset can yield further improvements in ROUGE performance. To check this, we created three sets of reference summaries from XSum articles using GPT-4, each varying in size, to serve as an enlarged training corpus. The specifics of the three datasets are detailed in Table 4.

First, we trained the model starting from the checkpoint fine-tuned on XSum, employing the same experimental setup as detailed in 3.4, results are reported in Table 5. On analysis, it becomes evident that augmenting the size of the dataset leads to an improvement in model performance, as measured by the ROUGE metric.

However, as we use the XSum checkpoint for its proven quality as a baseline, human reference remains crucial in our training process, leading to redundancy compared to other systems. To address this redundancy, we conducted additional experiments where, alongside using the XSum checkpoint, we initiated training with pre-trained weights for each model in this new configuration. Recognizing the influence of data volume on training efficacy, our performance evaluation in this experiment was confined to only two dataset sizes, medium and large.

| Variant | Train | Validation | Test |
|--------------|---------|------------|--------|
| Small (20k) | 20,000 | 1,100 | |
| Medium (50K) | 50,000 | 2,750 | 11,334 |
| Large (100K) | 100,000 | 5,500 | |

Table 4: Details of three dataset variations on XSum

| System | Dataset | R-1 | R-2 | R-L |
|------------------------|---------|--------------|--------------|--------------|
| BART _{GPT} | Small | 45.36 | 19.59 | 36.28 |
| | Medium | 47.44 | 21.47 | 38.34 |
| | Large | 48.52 | 22.42 | 39.57 |
| PEGASUS _{GPT} | Small | 43.71 | 18.68 | 35.07 |
| | Medium | 46.63 | 20.99 | 38.12 |
| | Large | 47.62 | 22.13 | 39.32 |
| BRIO _{GPT} | Small | 47.37 | 21.30 | 38.55 |
| | Medium | 48.82 | 23.28 | 40.66 |
| | Large | 49.05 | 23.81 | 41.20 |

Table 5: Evaluation of ROUGE scores after fine-tuning from the XSum checkpoint with various data sizes.

| System | Dataset | R-1 | R-2 | R-L |
|------------------------|---------|--------------|--------------|--------------|
| BART _{GPT} | Small | 46.28 | 20.37 | 37.26 |
| | Medium | 48.06 | 22.11 | 39.60 |
| | Large | 48.84 | 23.13 | 40.68 |
| PEGASUS _{GPT} | Small | 45.04 | 19.59 | 36.26 |
| | Medium | 47.21 | 21.76 | 38.80 |
| | Large | 47.88 | 22.50 | 39.64 |
| BRIO _{GPT} | Small | 47.64 | 21.68 | 38.93 |
| | Medium | 48.99 | 23.35 | 40.79 |
| | Large | 49.33 | 24.08 | 41.44 |

Table 6: Evaluation of ROUGE scores post fine-tuning from pre-trained weight with different data sizes.

Table 6 shows the result of fine-tuning from pre-trained weight. We observed that:

- (1) The model performance can indeed be advanced by training with only LLM reference, **which proved that our dataset can substitute the original XSum dataset in the training procedure.**
- (2) Compared to the model fine-tuned on the XSum checkpoint, the model that was fine-tuned from pre-trained weights demonstrated enhanced performance on identical data volumes. This improve-

ment likely originates from variances between human reference and LLM reference (detailed in section 4.2), prompting the model to perceive previously trained targets as potential noise.

(3) **Our dataset reduces the performance gap across models like BART, PEGASUS, and BRIO**, indicating that summaries generated using LLM effectively counteract biases associated with the varied styles of human writers in the original dataset. Therefore, these LLM-generated summaries facilitate a smoother learning process for models, thereby diminishing the requirement for intricate training methodologies.

5.2 Novelty Analysis

In this section, we delve into the comparative analysis of novelty between the summaries authored by humans and those generated by GPT-4. Novelty is defined through the computation of novel n-grams, a method that serves to gauge the 'abstraction' of our models. The novelty metric is calculated⁸ using the formula from Liu et al. (2022), i.e.,

$$Novelty(D, S^*) = \frac{\sum_{g \in G_{S^*}} \mathbb{1}(g \notin G_D)}{|G_{S^*}|} \quad (1)$$

where D and S^* are the source document and reference summary respectively, G_D and G_{S^*} are the sets of bigrams in D and S^* , $\mathbb{1}$ is the indicator function.

As presented in Table 7, models referencing GPT-4 exhibit better abstraction compared to those referencing human-generated summaries in the CNNDM dataset. Conversely, for the XSum dataset, models using human references are more "abstract" than those based on GPT-4 references. Despite these differences, as discussed in Section 4, summaries guided by GPT-4 are favored by human annotators across both the XSum and CNNDM datasets. This preference suggests that GPT-4, alongside our model, successfully balances the use of a diverse vocabulary for summary composition with effective information extraction from the original articles. Such a balance enhances summary relevance and aligns more closely with human preferences in summary generation.

6 Conclusion

In this work, we propose a novel supervised learning framework that leverages summaries generated

⁸The calculation is performed using ExplainaBoard (Liu et al., 2021). <https://github.com/neulab/ExplainaBoard>, and we had not employed PTBTokenizer prior to this calculation.

| System | XSum | | CNNDM | |
|------------------------|--------------|--------------|--------------|--------------|
| | Unigram | Bigram | Unigram | Bigram |
| Human | .3399 | .8342 | .1180 | .4960 |
| GPT-4 | .2960 | .8009 | .2375 | .7074 |
| BART | .2461 | .7310 | .0118 | .0922 |
| BART _{GPT} | .1986 | .6643 | .1287 | .5389 |
| PAGASUS | .2664 | .7474 | .1666 | .2919 |
| PAGASUS _{GPT} | .1558 | .5780 | .0946 | .4616 |
| BRIO | .2696 | .7654 | .0258 | .2261 |
| BRIO _{GPT} | .2203 | .7039 | .1389 | .5666 |

Table 7: Ratio of novel n -grams of various models on XSum and CNNDM. Novel n -grams are those that appear in the summaries but not in the source documents.

using LLMs as references. We performed an extensive human evaluation to compare systems guided by human-written summaries and those guided using LLM-generated summaries, analyzing the produced summaries across various dimensions. Our findings suggest that LLMs can guide small summarization models to produce summaries closely aligned with human preferences, indicating a new direction for research in automatic summarization. Furthermore, to facilitate ongoing research, we are releasing **GXSum** datasets in three sizes, comprising articles from the XSum dataset and summaries generated using LLMs. Our experiments validate the potential of our dataset to replace the original XSum dataset. We believe that our insights and the datasets we provide will encourage further exploration into the application of LLM knowledge in enhancing smaller, task-specific language models. We believe that our findings and released dataset provide new and unique insights into the LLM-enhanced automatic text summarization task.

7 Limitations

Our work introduces a new dataset, GXsum, for which we employ summaries generated by GPT-4 as references. It is essential to note that in our experiments, summaries were generated using OpenAI's API, which, due to its rapid iteration capability, might result in variable outcomes that could limit the reproducibility of our experiments. Furthermore, constrained by the performance of GPT-4, the generated summaries may still possess a certain level of hallucination. Additionally, considering effectiveness, the dataset and generated summaries used in this experiment are confined to the news domain. Employing datasets from other domains might provide a more comprehensive anal-

570 ysis, which represents a potential future research
 571 direction for us. Lastly, the human evaluation ex-
 572 periments conducted aim to explore a wide range
 573 of human reading preferences. The outcomes may
 574 vary depending on the timing of the assessment and
 575 the platform used to employ evaluators; we merely
 576 state the observed facts.

577 References

578 Toufique Ahmed and Premkumar Devanbu. 2023.
 579 [Few-shot training llms for project-specific code-](#)
 580 [summarization](#). In *Proceedings of the 37th*
 581 *IEEE/ACM International Conference on Automated*
 582 *Software Engineering*, ASE '22, New York, NY,
 583 USA. Association for Computing Machinery.

584 Lochan Basyal and Mihir Sanghvi. 2023. [Text summa-](#)
 585 [rization using large language models: A comparative](#)
 586 [study of mpt-7b-instruct, falcon-7b-instruct, and ope-](#)
 587 [nai chat-gpt models](#).

588 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
 589 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
 590 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
 591 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
 592 Gretchen Krueger, Tom Henighan, Rewon Child,
 593 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
 594 Clemens Winter, Christopher Hesse, Mark Chen, Eric
 595 Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,
 596 Jack Clark, Christopher Berner, Sam McCandlish,
 597 Alec Radford, Ilya Sutskever, and Dario Amodei.
 598 2020. Language models are few-shot learners. In
 599 *Proceedings of the 34th International Conference on*
 600 *Neural Information Processing Systems*, NIPS'20,
 601 Red Hook, NY, USA. Curran Associates Inc.

602 Chris Callison-Burch and Mark Dredze. 2010. Creating
 603 speech and language data with amazon’s mecha-
 604 nical turk. In *Proceedings of the NAACL HLT 2010*
 605 *workshop on creating speech and language data with*
 606 *Amazon’s Mechanical Turk*, pages 1–12.

607 Bosheng Ding, Chengwei Qin, Linlin Liu, Yew Ken
 608 Chia, Boyang Li, Shafiq Joty, and Lidong Bing. 2023.
 609 [Is GPT-3 a good data annotator?](#) In *Proceedings*
 610 *of the 61st Annual Meeting of the Association for*
 611 *Computational Linguistics (Volume 1: Long Papers)*,
 612 pages 11173–11195, Toronto, Canada. Association
 613 for Computational Linguistics.

614 Zi-Yi Dou, Pengfei Liu, Hiroaki Hayashi, Zhengbao
 615 Jiang, and Graham Neubig. 2021. [GSum: A gen-](#)
 616 [eral framework for guided neural abstractive summa-](#)
 617 [rization](#). In *Proceedings of the 2021 Conference of*
 618 *the North American Chapter of the Association for*
 619 *Computational Linguistics: Human Language Tech-*
 620 *nologies*, pages 4830–4842, Online. Association for
 621 Computational Linguistics.

622 Alexander R. Fabbri, Wojciech Kryściński, Bryan Mc-
 623 Cann, Caiming Xiong, Richard Socher, and Dragomir

Radev. 2021. [SummEval: Re-evaluating Summariza-](#)
 tion Evaluation. *Transactions of the Association for*
Computational Linguistics, 9:391–409. 624
625
626

Zorik Gekhman, Jonathan Herzig, Roei Aharoni, Chen
 Elkind, and Idan Szpektor. 2023. [TrueTeacher:](#)
[Learning factual consistency evaluation with large](#)
[language models](#). In *Proceedings of the 2023 Con-*
ference on Empirical Methods in Natural Language
Processing, pages 2053–2070, Singapore. Associa-
 tion for Computational Linguistics. 627
628
629
630
631
632
633

Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022.
 News summarization and evaluation in the era of
 gpt-3. *arXiv preprint arXiv:2209.12356*. 634
635
636

Cheng-Yu Hsieh, Chun-Liang Li, Chih-kuan Yeh,
 Hootan Nakhost, Yasuhisa Fujii, Alex Ratner, Ranjay
 Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. [Dis-](#)
[tilling step-by-step! outperforming larger language](#)
[models with less training data and smaller model](#)
[sizes](#). In *Findings of the Association for Compu-*
tational Linguistics: ACL 2023, pages 8003–8017,
 Toronto, Canada. Association for Computational Lin-
 guistics. 637
638
639
640
641
642
643
644
645

Junmo Kang, Wei Xu, and Alan Ritter. 2023. [Distill or](#)
[annotate? cost-efficient fine-tuning of compact mod-](#)
[els](#). In *Proceedings of the 61st Annual Meeting of the*
Association for Computational Linguistics (Volume 1:
Long Papers), pages 11100–11119, Toronto, Canada.
 Association for Computational Linguistics. 646
647
648
649
650
651

Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis,
 Jelena Luketina, Eric Hambro, Edward Grefenstette,
 and Roberta Raileanu. 2023. Understanding the ef-
 fects of rlhf on llm generalisation and diversity. *arXiv*
preprint arXiv:2310.06452. 652
653
654
655
656

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan
 Ghazvininejad, Abdelrahman Mohamed, Omer Levy,
 Veselin Stoyanov, and Luke Zettlemoyer. 2020.
[BART: Denoising sequence-to-sequence pre-training](#)
[for natural language generation, translation, and com-](#)
[prehension](#). In *Proceedings of the 58th Annual Meet-*
ing of the Association for Computational Linguistics,
 pages 7871–7880, Online. Association for Computa-
 tional Linguistics. 657
658
659
660
661
662
663
664
665

Shiyang Li, Jianshu Chen, Yelong Shen, Zhiyu Chen,
 Xinlu Zhang, Zekun Li, Hong Wang, Jing Qian,
 Baolin Peng, Yi Mao, Wenhui Chen, and Xifeng
 Yan. 2022. [Explanations from large language models](#)
[make small reasoners better](#). 666
667
668
669
670

Wei Li, Xinyan Xiao, Yajuan Lyu, and Yuanzhuo Wang.
 2018. [Improving neural abstractive document sum-](#)
[marization with explicit information selection mod-](#)
[eling](#). In *Proceedings of the 2018 Conference on*
Empirical Methods in Natural Language Processing,
 pages 1787–1796, Brussels, Belgium. Association
 for Computational Linguistics. 671
672
673
674
675
676
677

Chin-Yew Lin. 2004. Rouge: A package for automatic
 evaluation of summaries. In *Text summarization*
branches out, pages 74–81. 678
679
680

| | | |
|-----|---|-----|
| 681 | Pengfei Liu, Jinlan Fu, Yang Xiao, Weizhe Yuan, Shuaichen Chang, Junqi Dai, Yixin Liu, Zihuiwen Ye, and Graham Neubig. 2021. ExplainaBoard: An explainable leaderboard for NLP . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations</i> , pages 280–289, Online. Association for Computational Linguistics. | 739 |
| 682 | | 740 |
| 683 | | 741 |
| 684 | | 742 |
| 685 | | 743 |
| 686 | | |
| 687 | | 744 |
| 688 | | 745 |
| 689 | | 746 |
| 690 | Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruo Chen Xu, and Chenguang Zhu. 2023a. G-eval: NLG evaluation using gpt-4 with better human alignment . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2511–2522, Singapore. Association for Computational Linguistics. | 747 |
| 691 | | 748 |
| 692 | | 749 |
| 693 | | 750 |
| 694 | | |
| 695 | | 751 |
| 696 | | 752 |
| 697 | Yixin Liu, Alex Fabbri, Pengfei Liu, Yilun Zhao, Linyong Nan, Ruilin Han, Simeng Han, Shafiq Joty, Chien-Sheng Wu, Caiming Xiong, and Dragomir Radev. 2023b. Revisiting the gold standard: Grounding summarization evaluation with robust human evaluation . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4140–4170, Toronto, Canada. Association for Computational Linguistics. | 753 |
| 698 | | 754 |
| 699 | | 755 |
| 700 | | 756 |
| 701 | | |
| 702 | | 757 |
| 703 | | 758 |
| 704 | | 759 |
| 705 | | 760 |
| 706 | Yixin Liu, Pengfei Liu, Dragomir Radev, and Graham Neubig. 2022. BRIO: Bringing order to abstractive summarization . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2890–2903, Dublin, Ireland. Association for Computational Linguistics. | 761 |
| 707 | | 762 |
| 708 | | |
| 709 | | 763 |
| 710 | | 764 |
| 711 | | 765 |
| 712 | | 766 |
| 713 | Yixin Liu, Kejian Shi, Katherine S He, Longtian Ye, Alexander R. Fabbri, Pengfei Liu, Dragomir Radev, and Arman Cohan. 2023c. On learning to summarize with large language models as references . | 767 |
| 714 | | 768 |
| 715 | | 769 |
| 716 | | 770 |
| 717 | Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1906–1919, Online. Association for Computational Linguistics. | 771 |
| 718 | | 772 |
| 719 | | 773 |
| 720 | | |
| 721 | | 774 |
| 722 | | 775 |
| 723 | Ramesh Nallapati, Bowen Zhou, Cicero dos Santos, Çağlar Gulçehre, and Bing Xiang. 2016a. Abstractive text summarization using sequence-to-sequence RNNs and beyond . In <i>Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning</i> , pages 280–290, Berlin, Germany. Association for Computational Linguistics. | 776 |
| 724 | | 777 |
| 725 | | 778 |
| 726 | | 779 |
| 727 | | 780 |
| 728 | | 781 |
| 729 | | 782 |
| 730 | Ramesh Nallapati, Bowen Zhou, Cicero dos Santos, Çağlar Gulçehre, and Bing Xiang. 2016b. Abstractive text summarization using sequence-to-sequence RNNs and beyond . In <i>Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning</i> , pages 280–290, Berlin, Germany. Association for Computational Linguistics. | 783 |
| 731 | | 784 |
| 732 | | 785 |
| 733 | | |
| 734 | | 786 |
| 735 | | 787 |
| 736 | | 788 |
| 737 | Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don’t give me the details, just the summary! | 789 |
| 738 | | 790 |
| | | 791 |
| | | 792 |
| | | 793 |
| | | |
| | | 744 |
| | | 745 |
| | | 746 |
| | | |
| | | 747 |
| | | 748 |
| | | 749 |
| | | 750 |
| | | |
| | | 751 |
| | | 752 |
| | | 753 |
| | | 754 |
| | | 755 |
| | | 756 |
| | | |
| | | 757 |
| | | 758 |
| | | 759 |
| | | 760 |
| | | 761 |
| | | 762 |
| | | |
| | | 763 |
| | | 764 |
| | | 765 |
| | | 766 |
| | | 767 |
| | | |
| | | 768 |
| | | 769 |
| | | 770 |
| | | 771 |
| | | 772 |
| | | 773 |
| | | |
| | | 774 |
| | | 775 |
| | | 776 |
| | | 777 |
| | | 778 |
| | | 779 |
| | | |
| | | 780 |
| | | 781 |
| | | 782 |
| | | 783 |
| | | 784 |
| | | 785 |
| | | |
| | | 786 |
| | | 787 |
| | | 788 |
| | | 789 |
| | | 790 |
| | | 791 |
| | | 792 |
| | | 793 |

Yizhong Wang, Swaroop Mishra, Pegah Alipoor-molabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, et al. 2022b. Benchmarking generalization via in-context instructions on 1,600+ language tasks. *arXiv e-prints*, pages arXiv–2204.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. **Transformers: State-of-the-art natural language processing**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J. Liu. 2020. Pegasus: pre-training with extracted gap-sentences for abstractive summarization. In *Proceedings of the 37th International Conference on Machine Learning, ICML'20*. JMLR.org.

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*, 12:39–57.

A LLM Summary Generation

A.1 Prompt Template Example

Assuming you are an abstract writer, responsible for writing summaries of articles. Given the source article: {article}, please write a summary between {len_lower} to {len_upper} words about this article. please ensure that the summary is grammatically correct and coherent.

Figure 3: Template for a ChatGPT API prompt.

Figure 3 illustrates the template for our prompt design. The {article} variable represents the source article from the original dataset, and the {len_lower} and {len_upper} variables represent the lower bound and upper bound length constraints that we will set.

A.2 Generation Process

Figure 4 shows an example of our LLM summary generation process.

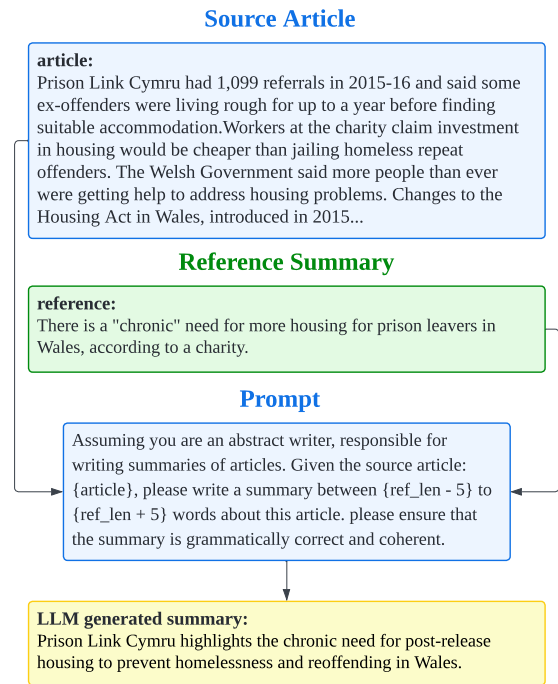


Figure 4: Illustration of LLM summary generation process

B Human Annotation Setting

B.1 Annotation Guideline

The definitions of various quality aspects we use in our annotation tasks are as follows:

- **Relevance:** Measures the importance of the summary content relative to the article, considering whether it has extracted the key points.
- **Consistency:** Considers whether the summary accurately includes all facts without fabricating false information.
- **Fluency:** Assesses whether each sentence in the summary is well-written and grammatically correct.
- **Coherence:** Considers whether the entire summary flows smoothly and reads naturally.
- **Informativeness:** Considers whether the summary clearly conveys the main message of the article, excluding unnecessary details.

B.2 Amazon Mechanical Turk Recruitment

To recruit qualified crowd annotators, stringent selection criteria were applied. These criteria stipulated that participants must have successfully completed more than 500 Human Intelligence Tasks (HITs), maintained an acceptance rate exceeding

858 95%, and resided within the United States. This rig-
859 orous selection process was implemented to guar-
860 antee that the annotators were native English speak-
861 ers and had a demonstrable record of experience in
862 effectively performing annotation tasks.

863 **B.3 Upwork Recruitment**

864 To ensure the rigor of expert analysis, we estab-
865 lished specific criteria for selecting annotators, fo-
866 cusing on those with a requisite level of expertise.
867 We engaged the Upwork platform to identify suit-
868 able candidates, stipulating prerequisites such as
869 residency in English-speaking countries (specifi-
870 cally the USA, UK, Australia, or Canada), min-
871 imum educational attainment of a bachelor’s de-
872 gree, and prior experience in data annotation or
873 linguistics-related roles. The ultimate selection of
874 our expert candidates comprised individuals with
875 backgrounds as writers, journalists, and profes-
876 sional text data annotators.

877 **C Case Study**

| Article ID | 36043765 | AVG. Score |
|------------------------------|---|------------|
| Human | A seal found tangled in nets on an Aberdeenshire beach has been returned to the sea. | 1.0 |
| GPT-4 | Scottish SPCA rescues and frees a heavily netted grey seal at Cruden Bay, preventing potential lethal injuries. | 3.0 |
| BART_{GPT} | A large seal entangled in netting at Cruden Bay was rescued by the Scottish SPCA and safely released. | 3.3 |
| PEGASUS_{GPT} | Scottish SPCA rescued a grey seal from Cruden Bay with a large net tangled around its neck, releasing him unharmed. | 3.0 |
| BRIO_{GPT} | A grey seal entangled in heavy netting at Cruden Bay was rescued by the Scottish SPCA and released unharmed. | 4.6 |

Table 8: Case study on XSum

| Article ID | 38537698 | AVG. Score |
|------------------------------|---|------------|
| Human | The reaction from BT’s investors told us much about media regulator Ofcom’s ruling on the fate of Openreach, the BT subsidiary that provides much of the UK’s broadband infrastructure. | 1.6 |
| GPT-4 | BT’s shares surged after Ofcom ruled out a company break-up due to practical challenges, including a complex pension scheme and legal obstacles, potentially leading to future operational conflicts. | 3.3 |
| BART_{GPT} | Ofcom’s Sharon White admits BT’s potential break-up faces practical hurdles due to land deals, pension scheme complexities, and potential conflicts over Openreach’s ownership and governance. | 4.3 |
| PEGASUS_{GPT} | BT shares rise 3% after Ofcom’s decision not to break up the company, citing pension issues and land deals. | 3 |
| BRIO_{GPT} | Ofcom postpones BT’s break-up due to pension scheme and land deals, prompting plans for separation and potential conflicts of interest. | 2.6 |

Table 9: Case study on XSum

| Article ID | ee17dfb574fec82ccac5689595e47483bd23f12 | AVG. Score |
|------------------------------|---|-------------------|
| Human | London's Metropolitan Police say the man was arrested at Luton airport after landing on a flight from Istanbul. He's been charged with terror offenses allegedly committed since the start of November. | 1.0 |
| GPT-4 | British man, Yahya Rashid, 19, faces terror charges upon his UK return from Turkey. Arrested at Luton airport, he's accused of preparing and aiding terrorist acts from November to March. Rashid will appear in court in Westminster. | 2.6 |
| BART_{GPT} | Yahya Rashid, a 19-year-old from northwest London, was charged with terror offenses after his arrest at Luton airport on his return from Turkey. He faces charges of terrorism preparation and aiding acts of terrorism between November 1 and March 31, with a court appearance set for Wednesday. | 5.0 |
| PEGASUS_{GPT} | Yahya Rashid, a 19-year-old from London, was charged with terrorism offenses at Luton Airport after returning from Turkey. He faces charges of preparing acts of terrorism and assisting others to commit terrorism. | 2.6 |
| BRIO_{GPT} | 19-year-old Yahya Rashid, a UK man, was charged with terror offenses after his arrest at London's Luton airport after his return from Turkey. He faces charges for planning and aiding acts of terrorism between November 1 and March 31, with his court appearance set for Wednesday. | 3.6 |

Table 10: Case Study on CNNDM.