

Are LLMs Better than Reported?

Detecting Label Errors and Mitigating Their Effect on Model Performance

Anonymous ACL submission

Abstract

NLP benchmarks rely on standardized datasets for training and evaluating models and are crucial for advancing the field. Traditionally, expert annotations ensure high-quality labels; however, the cost of expert annotation does not scale well with the growing demand for larger datasets required by modern models. While crowd-sourcing provides a more scalable solution, it often comes at the expense of annotation precision and consistency. Recent advancements in large language models (LLMs) offer new opportunities to enhance the annotation process, particularly for detecting label errors in existing datasets. In this work, we consider the recent approach of LLM-as-a-judge, leveraging an ensemble of LLMs to flag potentially mislabeled examples. Through a case study of four datasets from the TRUE benchmark, covering different tasks, we empirically analyze the labeling quality of existing datasets and compare expert, crowd-sourced, and LLM-based annotations in terms of the agreement, label quality, and efficiency, demonstrating the strengths and limitations of each annotation method. Our findings reveal a substantial number of label errors, which, when corrected, induce a significant upward shift in reported model performance. This suggests that many of the LLMs’ so-called mistakes are due to label errors rather than genuine model failures. Additionally, we discuss the implications of mislabeled data and propose methods to mitigate them in training to improve performance.

1 Introduction

Natural Language Processing (NLP) benchmarks have long served as a cornerstone for advancing the field, providing standardized datasets for training and evaluating methods and models (Wang et al., 2019; Hendrycks et al., 2021; Srivastava et al., 2023; Calderon et al., 2024). These datasets have been developed over the years for various tasks and scales, annotated using different schemes. Gold

labels represent the “true” or ground truth annotations, which are typically established through expensive rigorous processes, including expert consensus and extensive quality control. However, as models have increased in size (Devlin et al., 2019; Brown et al., 2020), the demand for larger datasets has also grown (Kaplan et al., 2020). Since expert annotation is cost-prohibitive, it does not scale well to meet these demands. The demand for large quantities of annotated data quickly and cost-effectively has led researchers to adopt crowd-sourcing, often sacrificing expertise for scale.

That way or another, constructing datasets heavily involves making compromises in annotation, trading off between scale, efficiency and expertise. Even when annotated by experts, datasets can naturally contain labeling errors, arising from factors such as task subjectivity, annotator fatigue, inattention, insufficient guidelines, and more (Rogers et al., 2013; Reiss et al., 2020; Syloypavan et al., 2023). Mislabeled data is even more pronounced when non-expert annotators are involved (Kennedy et al., 2020; Chong et al., 2022a). Widespread mislabeled data is particularly concerning because both the research community and the industry rely heavily on benchmarks. In training data, label errors harm model quality and hinder generalization, while in test sets, they lead to flawed comparisons, false conclusions, and prevent progress.

Recent advancements in large language models (LLMs) (Ouyang et al., 2022; Chiang and Lee, 2023; Li et al., 2023; Gat et al., 2024) present new opportunities to improve the annotation process, specifically in detecting label errors within existing datasets. Rather than re-annotating entire datasets (e.g., through experts or crowd-workers), we consider the LLM-as-a-judge approach (Zheng et al., 2023), and propose a simple yet effective method by leveraging an ensemble of LLMs to flag a set of potentially mislabeled examples. These can then be sent to experts for re-annotation and correction,

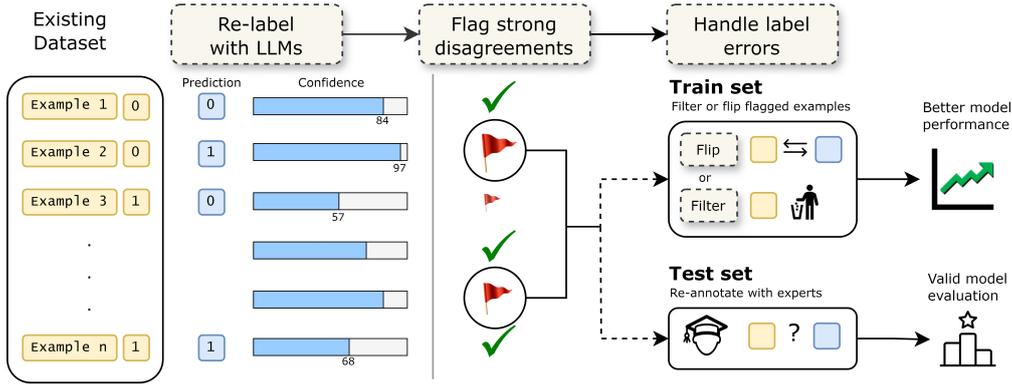


Figure 1: An illustration of our approach for detecting and addressing mislabeled data: (1) Re-label examples from existing datasets using an ensemble of LLMs. (2) Identify *strong disagreements* between the LLM’s predictions and the original labels (i.e., high confidence in a different label), flagging examples based on confidence levels. Our findings show that LLMs detect between 6% and 21% of label errors, and higher LLM confidence is strongly associated with improved precision in error detection. (3) In the training set, we either filter or flip flagged examples, leading to an increase of up to 4%. For the test set, flagged examples are re-annotated by experts to make sure the evaluation is accurate. Under accurate evaluation, the performance of LLMs is up to 15% higher.

or even get filtered during training.

Specifically, we construct an ensemble model using multiple LLMs with diverse prompts, gathering both their predicted labels and corresponding confidence scores. These predictions are contrasted with the original labels, and instances where the LLMs *strongly disagree* with the original label (i.e., show high confidence in a different label) are flagged as potential mislabeling cases. Additionally, we not only explore the role of LLMs in detecting errors but also evaluate their performance as annotators, comparing them with expert and crowd-sourced annotations. We assess these approaches in terms of agreement, label quality, and efficiency, highlighting their strengths and limitations.

We aim to answer the following questions through a comprehensive end-to-end study: (1) Do current benchmarks include mislabeled data? (2) Can LLMs detect label errors? (3) How do expert, crowd-sourced, and LLM-based annotations compare in quality and efficiency? and (4) What are the implications of mislabeled data on model performance and can we mitigate their impact?

To this end, we choose the TRUE benchmark (Honovich et al., 2022) – A collection consolidating 11 existing datasets annotated for factual consistency in a unified format – as a case-study and empirically investigate its labeling quality. Specifically, we analyze four datasets from TRUE with binary factual consistency annotation originating from different tasks. To support our claims and results in other setups, we conduct similar experi-

ments on an additional dataset, SummEval (Fabbri et al., 2021), which evaluates generated summaries in four dimensions on a scale of 1 to 5.

Our paper presents both methodological and empirical contributions. We propose a straightforward approach for detecting potential mislabeled examples (as illustrated in Figure 1), revealing a substantial number of label errors in existing datasets, ranging from 6% to 21%. Additionally, we demonstrate that the precision of LLMs in identifying errors improves with their confidence in an incorrect label; when their confidence exceeds 95%, over two-thirds of those labels are human errors. Moreover, we show that LLM-based annotations not only excel in error detection but also perform similarly to, or better than, traditional annotation methods, offering better trade-offs between quality, scale, and efficiency. Finally, we empirically illustrate the negative impact of mislabeled data on model training and evaluation. We propose a simple automated method for addressing label errors, improving the performance of fine-tuned models by up to 4%. In evaluation, we found that mislabeled data can significantly distort reported performance; LLMs may perform up to 15% better. This indicates that many so-called prediction errors are not genuine errors but are instead human annotation mistakes.

2 Related Work

Traditional Human Annotation Approaches

Crowdsourcing is widely used for annotating large-scale NLP datasets (Rajpurkar et al., 2016;

148 Williams et al., 2018; Wang et al., 2022), offering
149 rapid and scalable data collection. However, quality
150 control remains a challenge, with labeling inconsis-
151 tencies increasing as dataset complexity grows (Lu
152 et al., 2020; Allahbakhsh et al., 2013). Moreover,
153 as LLMs approach near-human performance (Chi-
154 ang and Lee, 2023; Chen and Ding, 2023), crowd
155 workers increasingly rely on these models for as-
156 sistance, further complicating annotation quality
157 (Veselovsky et al., 2023b,a). Expert annotation pro-
158 vides more reliable labels for domain-specific and
159 cognitively demanding tasks (e.g., medical or legal
160 domains) but is significantly slower and costlier
161 than crowdsourcing (Snow et al., 2008; Chau et al.,
162 2020). Ensuring inter-annotator agreement among
163 experts adds further complexity and expense (Bale-
164 dent et al., 2022). Hybrid approaches that com-
165 bine expert and crowd-sourced annotations help
166 balance cost and quality, though expert oversight
167 remains crucial for high-quality labels (Nguyen
168 et al., 2015). Our study compares expert, crowd-
169 sourced, and LLM-based annotation approaches in
170 terms of quality and efficiency.

171 **LLMs in the Annotation Loop** LLMs have
172 been increasingly utilized as annotators in various
173 NLP tasks, offering potential benefits in efficiency
174 and scalability. Several studies have demonstrated
175 that LLMs can effectively generate annotations
176 from scratch, sometimes outperforming human an-
177 notators or crowd workers (He et al., 2023; Gi-
178 lardi et al., 2023; Törnberg, 2023; Calderon and
179 Reichart, 2024). However, LLMs are not flaw-
180 less and cannot be considered gold-standard an-
181 notators when used alone. They may produce in-
182 correct annotations, especially in complex (Chen
183 et al., 2024), social (Ventura et al., 2023; Felkner
184 et al., 2024), emotional (Lissak et al., 2024), or low-
185 resource (Bhat and Varma, 2023) contexts. These
186 studies showed that LLMs can exhibit poor per-
187 formance and biases, highlighting the necessity
188 of human oversight to ensure quality or fairness.
189 To address this issue, several approaches for col-
190 laborative (Kim et al., 2024; Li et al., 2023) or
191 active learning (Zhang et al., 2023; Kholodna et al.,
192 2024) were suggested, where LLMs and humans
193 are both part of the annotation procedure. While
194 most research focuses on annotation from scratch,
195 our work employs an ensemble of LLMs to flag po-
196 tentially mislabeled data points in existing datasets.

197 **Handling Label Errors** Label errors (also re-
198 ferred to as label noise) in training and evaluation

199 datasets can significantly impair NLP model per-
200 formance and reliability (Fréney and Verleysen,
201 2014). Previous work mainly focuses on fine-tuned
202 models and typically identifies mislabeled exam-
203 ples based on the model’s low confidence or high
204 training loss (Chong et al., 2022b; Hao et al., 2020;
205 Pleiss et al., 2020; Northcutt et al., 2019). For ex-
206 ample, Chong et al. (2022b) showed that ranking
207 data points based on the training loss can help de-
208 tect errors. Once these high-loss or low-confidence
209 examples are flagged, they are typically filtered out
210 (Nguyen et al., 2019; Northcutt et al., 2019), cor-
211 rected automatically (Pleiss et al., 2020; Hao et al.,
212 2020), or re-labeled by human annotators (North-
213 cutt et al., 2021) to verify and improve dataset qual-
214 ity. Unlike previous works, we use an ensemble
215 of LLMs to flag only high-confidence false predic-
216 tions. Our results demonstrate that low-confidence
217 examples weakly correlates with errors, but high-
218 confidence in the false predictions strongly do.

219 3 LLM as an Annotator and Detector

220 This study aims to evaluate the potential of LLMs
221 in detecting mislabeled examples and compare
222 three annotation approaches: experts, crowdsourc-
223 ing, and LLMs. To this end, we use an ensem-
224 ble model that combines multiple LLMs with var-
225 ied prompts. The motivation for this ensemble is
226 twofold: first, we demonstrate that it enhances error
227 detection and aligns more closely with expert an-
228 notations while also decreases the variance; second,
229 it offers a simple approach that avoids the need for
230 complex model selection or extensive prompt engi-
231 neering, relying instead on the collective strength.

232 **Prediction and Confidence** To make a predic-
233 tion using the ensemble, we first extract class prob-
234 abilities of each LLM and prompt from the logits
235 of the representing class tokens (e.g., 0 or 1 for the
236 binary TRUE datasets, and 1 to 5 for the ordinal
237 SummEval). The probabilities are then normalized
238 to sum to 1. Next, we compute the average probab-
239 ility for each class across the ensemble and select the
240 class with the highest probability (argmax) as the
241 final prediction. The confidence in the prediction is
242 defined as the corresponding ensemble probability.
243 If the token probabilities are not accessible, they
244 can be approximated via sampling.

245 **Errors Detection** We re-label the dataset us-
246 ing the ensemble, keeping both the prediction and
247 confidence for each example. We then flag poten-
248 tially mislabeled examples where there is *strong*

disagreement between the ensemble prediction and the original label, specifically when the model exhibits high confidence in a false prediction. In the binary case, we examine only examples where the ensemble prediction differs from the original label. In the ordinal case, we examine examples where the difference between the original label and the ensemble prediction is strictly greater than 1 (e.g., 3 vs. 5, 1 vs. 5, 4 vs. 2, etc.). After examining these examples, only those with confidence exceeding a predefined threshold are flagged as potentially mislabeled. Our experiments show that as confidence in an incorrect prediction increases, the likelihood of the example being mislabeled also rises.

For test sets, flagged examples can be re-examined by human experts to verify their true labels. For training sets, the same approach can be applied, but we also propose an automated method to improve model training: flagged examples can either be removed from the dataset or have their labels corrected based on the ensemble prediction.

4 Experimental Setup

4.1 Data

As a case-study, we choose to explore the extensive and widely used TRUE benchmark (Honovich et al., 2022), which is typically used as an evaluation set (Steen et al., 2023; Gekhman et al., 2023; Wang et al., 2024; Zha et al., 2023). It consists of 11 datasets from various NLP tasks such as summarization and knowledge-grounded dialogue. This benchmark is unique in its approach of bringing multiple datasets and tasks into a unified schema of binary factual consistency labels. Each dataset is transformed from its original structure (e.g., a source document and a summary) into two input texts, *Grounding* and *Generated Text*, and a binary label indicating whether the generated text is factually consistent w.r.t the grounding. This enables us to examine multiple tasks and domains under the same umbrella at once while maintaining a unified binary-label schema. Specifically, we focus on four TRUE datasets, one from each task: MNBM – summarization evaluation (Maynez et al., 2020); BEGIN – grounded dialogue evaluation (Dziri et al., 2022); VitaminC – fact verification (Schuster et al., 2021); and PAWS – paraphrasing evaluation (Zhang et al., 2019). See Appendix E for additional details on these datasets.

For each of the four datasets, we randomly sampled 1000 examples (or the whole dataset if the

number of examples is smaller than 1000). These examples are annotated by LLMs. We set an evaluation (i.e., test set) based on 160 randomly sampled examples from each dataset (a total of 640), while the rest remain for training and validation (they will be relevant for subsection 7.1). In addition to the LLM annotations, the evaluation set is also re-annotated by two experts three crowd worker.

SummEval In addition to the TRUE benchmark, we replicate some of the experiments on the full SummEval benchmark (Fabbri et al., 2021). This benchmark includes 1600 generated summaries evaluated on *four dimensions* (relevance, fluency, coherence, consistency) by crowd-workers and experts. The SummEval benchmark is widely used for benchmarking reference-free automatic evaluation methods such as LLM-as-a-judge. In contrast to TRUE, the labeling scheme is *ordinal* on a scale of 1 to 5. For further information on the SummEval data and experimental setting, see Appendix A. Noteworthy, when researchers employ the SummEval benchmark, they use solely the expert annotations. Accordingly, the focus of our experiments conducted on SummEval is (1) to simulate a setup where the original labels are obtained through crowd-sourcing while relying on expert annotations as the gold standard; and (2) to compare the three annotation approaches (crowd-sourcing, experts, and LLMs).

4.2 Annotation Procedure

This subsection outlines the annotation procedures for the various approaches. Refer to Appendix D for additional implementation and technical details not covered here, or Appendix A for the SummEval LLM annotation details.

LLMs We re-annotate the data with four LLMs: GPT-4, (OpenAI, 2023), PaLM2 (Anil et al., 2023), Mistral (7B) (Jiang et al., 2023), Llama 3 (8B) (Dubey et al., 2024), and GPT-4o and Gemini-1.5-Flash for SummEval. Our ensemble model leverages four different prompts which control the variance caused by task descriptions. The prompts are designed as a zero-shot classification task, e.g., for TRUE the requested output is a single token, either '0' for factual inconsistency or '1' for factual consistency (as described in Figure 12).

Crowd-sourcing Generally, crowd-sourced annotators span a spectrum– from untrained, "common" crowd-workers to carefully selected and trained annotators. Our paper focuses on the lower

end of this spectrum. We utilize the platform of Amazon Mechanical Turk (MTurk) to recruit crowd-workers for annotating 100 examples from each TRUE dataset (a total of 400), and to design the interface layout. Examples were randomly assigned to annotators. Each annotated example was manually reviewed. Rejected examples were returned to the pool and re-annotated until each example was annotated by three different annotators. To prevent LLM use, we disabled right-click and `Ctrl+C` in the platform (as suggested by Veselovsky et al., 2023a). To obtain a single label per example, we consider two different aggregations: (1) *Majority* - by majority vote, and (2) *Strict* - if any annotator marks it *inconsistent*, that becomes the label. For SummEval, we use the crowd-sourced annotations provided by Fabbri et al. (2021), aggregated by their median.

Experts All TRUE examples where the prediction differed from the original label, regardless of confidence, were annotated by human experts. The experts are two of the paper’s authors, who are fully familiar with the guidelines and task characteristics. Each example was independently annotated by both experts on a scale from 0 (*inconsistent*) to 1 (*consistent*). The examples were shuffled and presented in no specific order, with neither the original nor LLM labels shown. For cases where the experts disagreed, a reconciliation phase followed, during which they discussed and attempted to resolve their differences. For more details on the procedure and annotation platform, see Appendix D.2. After re-annotating all conflicted examples, we define the *gold label* as the original label, if the LLM prediction agrees with it, or the expert resolution, if there was a disagreement. For SummEval, we use the expert annotations provided by Fabbri et al. (2021), aggregated by their median.

5 Label Errors: Analysis and Detection

5.1 Do current benchmarks include mislabeled data?

To address the first research question, we annotate the test-set of TRUE (as described in section 4 using LLMs. We then contrast these annotations with the original labels, to find disagreements. As shown in Table 2, the disagreement rate is significant and can be up to $\sim 40\%$ of the examples. An example of such disagreement is presented in Table 1. While this would typically suggest that the LLMs performed poorly, we chose to further

Dataset: BEGIN

Grounding: Hillary Clinton, the nominee of the Democratic Party for president of the United States in 2016, has taken positions on political issues while serving as First Lady of Arkansas (1979–81; 1983–92), First Lady of the United States (1993–2001);
Generated Text: She is the nominee in 2016.

Original Label: 0 **LLM p :** 0.98 **Gold Label:** 1

Explanation: She (Hillary Clinton) is indeed the nominee in 2016 as specifically stated in the grounding.

Table 1: Example of an annotation error in the original datasets, discovered by LLMs and corrected by experts. In Appendix Table 6 we provide additional examples.

Dataset	Task	% pos	% LLM disagree	% error
MNBM	Summarization	10.6	39.4	16.9 (11.6)
BEGIN	Dialogue	38.7	34.4	21.2 (15.8)
VitaminC	Fact Verification	52.5	17.5	8.1 (4.4)
PAWS	Paraphrasing	44.3	22.5	6.2 (3.0)

Table 2: Summary of LLM disagreement and label error rates across different datasets. %pos is the percentage of positive (i.e., the *consistent* class) examples in the data. % LLM disagree refers to the percentage of examples where the LLM label differs from the original one. % error indicates the error rate in the sampled test set, while the number in parentheses denotes the estimated lower bound of the error rate for the entire dataset.

investigate these cases and resolve the disagreements. To this end, we asked human experts to re-annotate the examples, allowing us to determine which is more accurate: the original label or the LLMs’ prediction.

Our findings show a considerable number of label errors for all examined datasets (see the %error column in Table 2). Based on the experts *gold label* and the sample sizes, we also estimate a lower bound for the total percentage of label errors in the full datasets. We employed the Clopper-Pearson exact method (Clopper and Pearson, 1934) to construct a 95% confidence interval for the binomial proportion, adjusted by a finite population correction (FPC) (see more details in Appendix G.1). We provide the lower bound of these confidence intervals in parentheses in Table 2, under the %error column. The lower bounds range from 3% in the PAWS dataset to 15.8% in the BEGIN dataset.

5.2 Can LLMs Detect Label Errors?

As described in subsection 5.1, we utilize LLMs to flag candidates for mislabeling, and indeed find label errors. In this subsection, we focus on the LLM viewpoint, exploring the effect of LLM confidence,

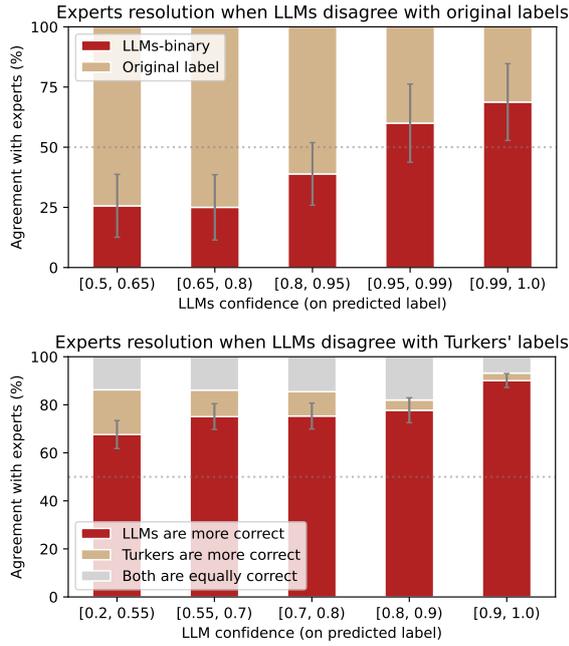


Figure 2: When LLMs disagree with original labels - who is correct? (**Top**) TRUE (**Bottom**) SummEval. As the LLM’s confidence grows, so does the precision of identifying an error in the original labels.

and the power of ensemble.

Confidence LLM annotations are valuable for flagging mislabeled data, offering more than just hard labels. By considering LLM confidence scores alongside their predictions, we can improve the precision of automatic error detection. Leveraging confidence can reduce re-annotation efforts by flagging only cases exceeding a predefined threshold. The rationale is that not all flagged examples should be treated equally. Instances flagged with low confidence indicate that the LLM recognizes a potential issue, however, when the LLM is highly confident in a label that contradicts the original one, it provides a stronger signal of a possible error.

Figure 2 shows the rate of the experts’ agreement with the LLMs compared to the agreement with original labels, divided into confidence-based bins. Bins are balanced by size, and defined by a confidence interval of 95% based on bootstrap sampling (see Appendix G.2 for further details). The bins reflect increasing levels of LLM confidence in its predicted label (i.e., a stronger disagreement between LLMs and the original labels).

From the top of Figure 2, we observe a clear trend: as LLM confidence increases, so does its precision in detecting label errors in the original dataset. In the highest confidence bin, LLM annotations surpass the original labels in agreement with

expert re-labeling, and this difference is statistically significant. This indicates that when the LLM is highly confident in its disagreement with the original label, the labeled example serves as a strong candidate for a labeling error. Note that even in cases where the expert agreement with LLMs was below 50%, mislabeled data was still discovered.

We replicated this analysis on the SummEval dataset (bottom of Figure 2) and observed a similar trend: higher confidence increases the likelihood that the LLM prediction is closer to the expert annotation than the original label. In the SummEval case, we consider the crowd-sourced labels as the original labels. For more details see Appendix A.

Ensemble By varying the size of the LLM ensemble, we examine two key aspects: predictive power (how well predictions align with gold labels, measured by ROC AUC for TRUE and average correlation for SummEval), and error detection power (measured by F1-score, averaging the recall of errors and the precision of correctly identifying a candidate as a true error). The ensemble power analysis is presented in Figure 3, with additional details in Appendix B. Our findings show that incorporating multiple LLMs and prompts in an ensemble is valuable. As the ensemble size increases, both label quality and error detection improve.

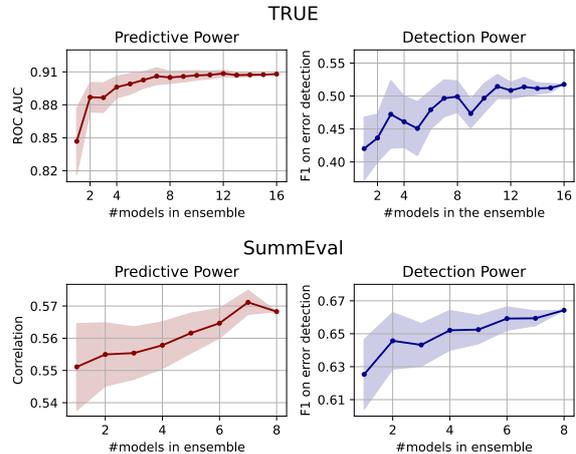


Figure 3: The power of ensemble. (**Top**) TRUE (**Bottom**) SummEval. As the ensemble size increases (x-axis), its performance against gold labels (Left), and its ability to detect label errors (Right) improves.

6 Comparing Annotation Approaches

Our paper discusses three annotation approaches, each with its own benefits and drawbacks, differing in how they balance label quality, scalability, and cost. Due to space limitations, we provide a

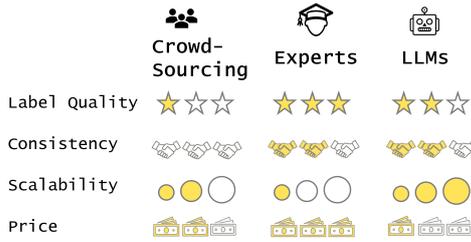


Figure 4: Annotation approaches comparison.

483 concise summary of our key findings here, with
 484 the full analysis available in Appendix C. Figure 4
 485 highlights the main insights.

486 **LLMs exhibit strong agreement with experts**
 487 **and among themselves.** Inter-annotator agreement
 488 (IAA) among LLMs, as well as their alignment
 489 with expert annotations, are significantly higher
 490 than that of crowd workers. In contrast, crowd-
 491 sourced annotations exhibit larger variability and
 492 lower agreement with experts, making them less
 493 reliable without additional verification.

494 **Crowd worker quality improves with experience**
 495 **but remains inconsistent.** Our analysis shows that
 496 experienced crowd workers produce higher-quality
 497 annotations. However, even among them, anno-
 498 tation quality and consistency remain lower than
 499 LLM-based annotation, which is more reliable.

500 **LLMs provide fast, scalable, and cost-efficient**
 501 **annotation.** Compared to expert and crowd-
 502 sourced annotation, LLMs require less time and are
 503 much more cost-effective per annotation, making
 504 them a viable alternative for large-scale annotation
 505 while effectively balancing the trade-off.

506 7 Implications of Mislabeled Data

507 7.1 Training on Mislabeled Data

508 Training on mislabeled data can harm model perfor-
 509 mance and stability, as learning from errors makes
 510 it harder to identify consistent patterns. The impact
 511 depends on various factors, such as the fraction
 512 of mislabeled data and the training procedure. In
 513 this subsection, we show that addressing this is-
 514 sue, even heuristically, significantly improves the
 515 model’s performance on a test set.

516 **Handling Label Errors** In order to handle label
 517 errors in the training set, and reduce its effect on
 518 model performance, we propose two manipulations.
 519 For both manipulations, we flag examples where
 520 the model strongly disagrees with the original label
 521 (i.e., with confidence above a certain threshold).

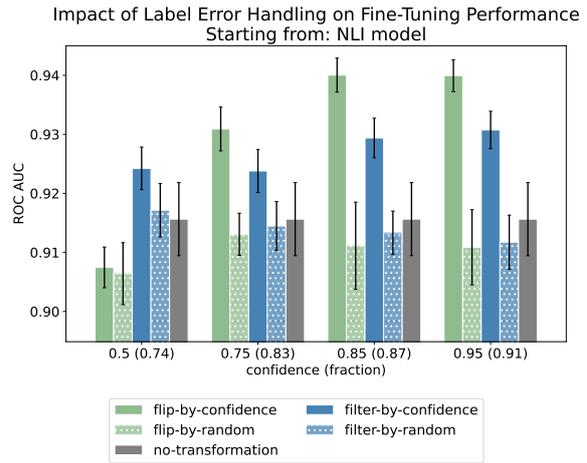


Figure 5: Fine-tuning a model on a transformed dataset. The gray bar is the original dataset - without any changes. The green bars present results for label flipping for a subset of examples, determined by LLMs-confidence (plain), or at random (dotted). The blue bars represent filtering of these examples.

522 The first manipulation is *filtering* flagged examples
 523 out, which maintains a “cleaner” yet smaller train-
 524 ing set. The second manipulation is label *flipping*
 525 for flagged examples, which maintains the same
 526 amount of data, but may also cause harm if flipping
 527 too many correct labels.

528 **Experimental Setup** We set the training set to
 529 be the additional data examples from the datasets
 530 (i.e., MNBM, BEGIN, VitaminC, PAWS), which
 531 are disjoint from the test set. Note that we posses
 532 gold labels for the test set alone, while for the train-
 533 ing set we only extract the confidence. The fine-
 534 tuning procedure includes splitting the training set
 535 into train and validation sets, and fine-tuning on the
 536 train set. We report average results of five seeds.

537 As an ablation study, we also apply these mani-
 538 pulations on a random subset of examples rather
 539 than the flagged examples. The ablation study aims
 540 to maintain a consistent number of training exam-
 541 ples, while the ablation for flipping aims to address
 542 the claim that in some cases, a relatively small
 543 fraction of label errors may be even considered as
 544 a noise that improves model robustness (e.g., as
 545 in label perturbation (Zhang et al., 2018) or label
 546 smoothing (Szegedy et al., 2016)).

547 We conducted this experiment starting from two
 548 base models: DeBERTa-v3, and a fine-tuned ver-
 549 sion of it on classic NLI datasets, which we will
 550 refer to as the NLI-base model. We chose the
 551 NLI-base model as NLI tasks closely resemble
 552 factual consistency evaluation (FCE), making it

Model	Rank		ROC AUC		F1 Score		Accuracy	
	Original	Gold	Original	Gold	Original	Gold	Original	Gold
GPT-4	3	1 (+2)	0.81	0.93 (+15%)	0.73	0.83 (+14%)	0.73	0.83 (+14%)
NLI model	1	2 (-1)	0.93	0.91 (-2%)	0.87	0.87 (—)	0.87	0.87 (—)
PaLM2	6	3 (+3)	0.81	0.91 (+12%)	0.71	0.81 (+14%)	0.71	0.81 (+14%)
GPT-4o	4	4 (—)	0.81	0.91 (+12%)	0.74	0.83 (+12%)	0.74	0.83 (+12%)
GPT-4-mini	5	5 (—)	0.81	0.91 (+12%)	0.71	0.79 (+11%)	0.70	0.79 (+13%)
Llama3	7	6 (+1)	0.75	0.86 (+15%)	0.47	0.50 (+6%)	0.52	0.55 (+6%)
Mistral-v0.3	8	7 (+1)	0.75	0.85 (+13%)	0.61	0.68 (+11%)	0.62	0.68 (+10%)
DeBERTa-v3	2	8 (-6)	0.84	0.80 (-5%)	0.76	0.73 (-4%)	0.76	0.73 (-4%)
Mistral-v0.2	9	9 (—)	0.73	0.82 (+12%)	0.66	0.72 (+9%)	0.66	0.72 (+9%)

Table 3: Comparison of Model Performance on Original and Gold Labels. Ranking is defined over ROC AUC.

well-suited for this experiment. Given the similar trends, we present the results for the NLI model here. Additional experiments and implementation details can be found in Appendix F.1.

Results Figure 5 shows the results of our experiments. In our confidence-based approaches, we clearly see the trend that as the confidence threshold—according to which our manipulations are applied—grows, our manipulation results in improved ROC AUC for both models. This trend eventually (i.e., for high enough LLM confidence) brings these approaches to significantly outperform the baseline. In contrast, when we applied our manipulations on random subsets, we generally see a diminishing effect of manipulation, converging to the no-manipulation baseline.

Comparing between the handling approaches, it appears that flipping is better than filtering for high confidence. We hypothesize that this stems from the amount of data that remains after flipping (i.e., the same amount as before the flipping) compared to the filtering approach, combined with the high error rate in these datasets. Note that this is contrary to the random case where filtering is better than flipping, as flipping a subset with low error-rate brings more damage than value.

7.2 Evaluating on Mislabeled Data

In this subsection, we examine the impact of mislabeled data in evaluation sets and its potential to distort results. Labeling errors can mislead the evaluation process, resulting in inaccurate performance metrics and, in some cases, flawed model comparisons that lead to incorrect conclusions.

Experimental Setup To test this assumption, we evaluate the performance of nine models, mostly state-of-the-art LLMs, on the test datasets. We compare their performance between the *original* labels, and the *gold* labels. For LLMs, we used zero-shot prediction as described in section 3, and averaged

over prompts. For DeBERTa-based models, we used the fine-tuned models from subsection 7.1, and averaged over seeds.

Results Prior to this work, an evaluation of these models would induce the values and ranking as in Table 3 under the *Original* sub-columns. However, as shown before, these datasets include labeling errors, and therefore do not support fair evaluation. Considering the new gold labels, based on expert intervention (as described in subsection 4.2), we obtain different results, shown in the *Gold* sub-columns. The first observed discrepancy is the ranking of models. For example, DeBERTa-v3 has shifted from being the second-best to the second-worst. Beyond the change in ranking, all metrics’ (i.e., ROC AUC, F1-score, and accuracy) range has shifted upward, indicating that LLMs perform better on this task than what was previously thought, likely due to label errors. If this phenomenon extends to other tasks and datasets beyond those examined in this study, it could suggest that LLMs are better than currently perceived.

8 Discussion

Labeling errors are a persistent issue in NLP datasets, negatively affecting model fine-tuning and evaluation. Our findings demonstrate that LLMs, particularly when highly confident, can effectively detect these errors, outperforming crowd workers in accuracy, consistency, and cost-efficiency. As LLM capabilities advance, their role in refining data quality will become central to improving NLP benchmarks. Future work could explore applying LLM-based error detection to a broader range of datasets and tasks, as well as refining methods for optimizing label correction strategies. We encourage researchers to adopt our methods and critically evaluate existing datasets to drive more robust, reliable results in the field.

630 Limitations

631 While our study provides valuable insights into the
632 role of LLMs in identifying label errors and im-
633 proving dataset quality, several limitations should
634 be considered. First, crowd workers encompass a
635 broad range of annotators with varying expertise
636 and training. Our analysis, focuses on the “com-
637 mon” crowd worker, typically an annotator selected
638 with minimal qualifications, such as an approved
639 task completion rate, and without specialized train-
640 ing. However, some datasets implement additional
641 measures, such as requiring prior experience or
642 task-specific instruction, which can influence anno-
643 tation quality. Importantly, we did not take crowd-
644 worker annotations at face value; we applied filter-
645 ing (based on the explanation crowd workers were
646 asked to write for each example) to remove a sub-
647 stantial number of low-quality assignments, such as
648 clearly invalid responses, in addition to enforcing
649 minimal qualification criteria.

650 Second, our analysis does not account for po-
651 tential data contamination, where LLMs may have
652 been trained on the datasets we evaluate. However,
653 since our analysis focuses on identifying and cor-
654 recting label errors within these datasets, contami-
655 nation would likely hinder rather than enhance our
656 findings. If an LLM had memorized these datasets,
657 it would be more likely to reproduce existing errors
658 rather than detect and correct them, making con-
659 tamination a potential limitation only for certain
660 aspects of evaluation but not for our core claims.

661 Third, LLM-based annotations can vary depend-
662 ing on the choice of prompting strategies and en-
663 semble methods. In this work, we use zero-shot
664 prompting and simple averaging for ensembling.
665 Still, alternative approaches – such as few-shot
666 prompting, chain-of-thought reasoning (Wei et al.,
667 2022), or self-refine (Madaan et al., 2023) – could
668 improve annotation accuracy and consistency. Like-
669 wise, for ensembling, more advanced methods—
670 such as percentile-based aggregation (Sherratt et al.,
671 2023), error-aware weighting (Freund and Schapire,
672 1997), confidence-aware methods (Lee, 2010; Lu
673 et al., 2024), or even LLM-based aggregation strate-
674 gies like debate variants (Liang et al., 2023; Du
675 et al., 2024) – may yield more reliable consensus
676 labels. We leave the exploration of these strate-
677 gies for future work and hope our study encourages
678 such further research.

References

- Mohammad Allahbakhsh, Boualem Benatallah, Alek- 680
sandar Ignjatovic, Hamid Reza Motahari-Nezhad, 681
Elisa Bertino, and Schahram Dustdar. 2013. [Quality 682](#)
[control in crowdsourcing systems: Issues and direc- 683](#)
[tions](#). *IEEE Internet Computing*, 17(2):76–81. 684
- Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin John- 685
son, Dmitry Lepikhin, Alexandre Passos, Siamak 686
Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng 687
Chen, Eric Chu, Jonathan H. Clark, Laurent El 688
Shafey, Yanping Huang, Kathy Meier-Hellstern, Gau- 689
rav Mishra, Erica Moreira, Mark Omernick, Kevin 690
Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, 691
Yuanzhong Xu, Yujing Zhang, Gustavo Hernández 692
Ábrego, Junwhan Ahn, Jacob Austin, Paul Barham, 693
Jan A. Botha, James Bradbury, Siddhartha Brahma, 694
Kevin Brooks, Michele Catasta, Yong Cheng, Colin 695
Cherry, Christopher A. Choquette-Choo, Aakanksha 696
Chowdhery, Clément Crepy, Shachi Dave, Mostafa 697
Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz, 698
Nan Du, Ethan Dyer, Vladimir Feinberg, Fangxi- 699
aoyu Feng, Vlad Fienber, Markus Freitag, Xavier 700
Garcia, Sebastian Gehrmann, Lucas Gonzalez, and 701
et al. 2023. [Palm 2 technical report](#). *CoRR*, 702
abs/2305.10403. 703
- Anaëlle Baledent, Yann Mathet, Antoine Widlöcher, 704
Christophe Couronne, and Jean-Luc Manguin. 2022. 705
[Validity, agreement, consensuality and annotated data 706](#)
[quality](#). In *International Conference on Language 707*
Resources and Evaluation. 708
- Savita Bhat and Vasudeva Varma. 2023. [Large lan- 709](#)
[guage models as annotators: A preliminary evalua- 710](#)
[tion for annotating low-resource language content](#). In 711
Proceedings of the 4th Workshop on Evaluation and 712
Comparison of NLP Systems, pages 100–107, Bali, 713
Indonesia. Association for Computational Linguis- 714
tics. 715
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie 716
Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind 717
Neelakantan, Pranav Shyam, Girish Sastry, Amanda 718
Askell, Sandhini Agarwal, Ariel Herbert-Voss, 719
Gretchen Krueger, Tom Henighan, Rewon Child, 720
Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, 721
Clemens Winter, Christopher Hesse, Mark Chen, Eric 722
Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, 723
Jack Clark, Christopher Berner, Sam McCandlish, 724
Alec Radford, Ilya Sutskever, and Dario Amodei. 725
2020. [Language models are few-shot learners](#). In *Ad- 726*
vances in Neural Information Processing Systems 33: 727
Annual Conference on Neural Information Process- 728
ing Systems 2020, NeurIPS 2020, December 6-12, 729
2020, virtual. 730
- Nitay Calderon, Naveh Porat, Eyal Ben-David, Alexan- 731
der Chapanin, Zorik Gekhman, Nadav Oved, Vitaly 732
Shalumov, and Roi Reichart. 2024. [Measuring the 733](#)
[robustness of nlp models to domain shifts](#). *arXiv 734*
preprint arXiv:2306.00168. 735

844	Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli.	2020. The shape of and solutions to the mturk quality crisis . <i>Political Science Research and Methods</i> , 8(4):614–629.	901
845	2023. Chatgpt outperforms crowd workers for text-annotation tasks . <i>Proceedings of the National Academy of Sciences of the United States of America</i> , 120.		902
846			903
847			
848			
849	Degan Hao, Lei Zhang, Jules H. Sumkin, Aly A. Mohamed, and Shandong Wu. 2020. Inaccurate labels in weakly-supervised deep learning: Automatic identification and correction and their impact on classification performance . <i>IEEE Journal of Biomedical and Health Informatics</i> , 24:2701–2710.	Nataliia Kholodna, Sahib Julka, Mohammad Khodadadi, Muhammed Nurullah Gumus, and Michael Granitzer. 2024. Llms in the loop: Leveraging large language model annotations for active learning in low-resource languages . <i>ArXiv</i> , abs/2404.02261.	904
850			905
851			906
852			907
853			908
854			
855	David N. Hauser, Aaron J. Moss, Cheskie Rosenzweig, Shalom N. Jaffe, Jonathan Robinson, and Leib Litman. 2021. Evaluating cloudresearch’s approved group as a solution for problematic data quality on mturk . <i>Behavior Research Methods</i> , 55:3953 – 3964.	Han Jun Kim, Kushan Mitra, Rafael Li Chen, Sajjadur Rahman, and Dan Zhang. 2024. Meganno+: A human-llm collaborative annotation system . In <i>Conference of the European Chapter of the Association for Computational Linguistics</i> .	909
856			910
857			911
858			912
859			913
860	Xingwei He, Zheng-Wen Lin, Yeyun Gong, Alex Jin, Hang Zhang, Chen Lin, Jian Jiao, Siu Ming Yiu, Nan Duan, and Weizhu Chen. 2023. Annollm: Making large language models to be better crowdsourced annotators . In <i>North American Chapter of the Association for Computational Linguistics</i> .	Klaus Krippendorff. 1970. Estimating the reliability, systematic error, and random error of interval data . <i>Educational and Psychological Measurement</i> , 30(1):61–70.	914
861			915
862			916
863			917
864			
865			
866	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.	Chi-Hoon Lee. 2010. Learning to combine discriminative classifiers: confidence based . In <i>Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Washington, DC, USA, July 25-28, 2010</i> , pages 743–752. ACM.	918
867			919
868			920
869			921
870			922
871			
872	Or Honovich, Roei Aharoni, Jonathan Herzig, Hagai Taitelbaum, Doron Kukliansky, Vered Cohen, Thomas Scialom, Idan Szpektor, Avinatan Hassidim, and Yossi Matias. 2022. TRUE: re-evaluating factual consistency evaluation . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022</i> , pages 3905–3920. Association for Computational Linguistics.	Minzhi Li, Taiwei Shi, Caleb Ziemis, Min-Yen Kan, Nancy F. Chen, Zhengyuan Liu, and Diyi Yang. 2023. Coannotating: Uncertainty-guided work allocation between human and large language models for data annotation . <i>ArXiv</i> , abs/2310.15638.	923
873			924
874			925
875			926
876			927
877			
878			
879			
880			
881			
882	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L�el�io Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth�ee Lacroix, and William El Sayed. 2023. Mistral 7b . <i>CoRR</i> , abs/2310.06825.	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujtu Yang, Zhaopeng Tu, and Shuming Shi. 2023. Encouraging divergent thinking in large language models through multi-agent debate . <i>CoRR</i> , abs/2305.19118.	928
883			929
884			930
885			931
886			932
887			
888			
889			
890	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 . <i>CoRR</i> , abs/2001.08361.	Shir Lissak, Nitay Calderon, Geva Shenkman, Yaakov Ophir, Eyal Fruchter, Anat Brunstein Klomek, and Roi Reichart. 2024. The colorful future of llms: Evaluating and improving llms as emotional supporters for queer youth . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), NAACL 2024, Mexico City, Mexico, June 16-21, 2024</i> , pages 2040–2079. Association for Computational Linguistics.	933
891			934
892			935
893			936
894			937
895			938
896			939
897			940
898			941
899			942
900			943
		Jian Lu, Wei Li, Qingren Wang, and Yiwen Zhang. 2020. Research on data quality control of crowdsourcing annotation: A survey . In <i>2020 IEEE Intl Conf on Dependable, Autonomic and Secure Computing, Intl Conf on Pervasive Intelligence and Computing, Intl Conf on Cloud and Big Data Computing, Intl Conf on Cyber Science and Technology Congress (DASC/PiCom/CBDCOM/CyberSciTech)</i> , pages 201–208.	944
			945
			946
			947
			948
			949
			950
			951
			952
		Zhihe Lu, Jiawang Bai, Xin Li, Zeyu Xiao, and Xinchao Wang. 2024. Beyond sole strength: Customized ensembles for generalized vision-language models . In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	953
			954
			955
			956
			957
			958

959	Bill MacCartney and Christopher D. Manning. 2009.	2022. Training language models to follow instructions with human feedback . In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .	1015
960	An extended model of natural logic . In <i>Proceedings of the Eight International Conference on Computational Semantics</i> , pages 140–156, Tilburg, The Netherlands. Association for Computational Linguistics.		1016
961			1017
962			1018
963			1019
964			1020
965	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhunoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: Iterative refinement with self-feedback . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	Geoff Pleiss, Tianyi Zhang, Ethan R. Elenberg, and Kilian Q. Weinberger. 2020. Identifying mislabeled data using the area under the margin ranking . <i>ArXiv</i> , abs/2001.10528.	1021
966			1022
967			1023
968			1024
969		Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.	1025
970			1026
971			1027
972			1028
973			1029
974			1030
975			
976	Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan T. McDonald. 2020. On faithfulness and factuality in abstractive summarization . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020</i> , pages 1906–1919. Association for Computational Linguistics.	Frederick Reiss, Hong Xu, Bryan Cutler, Karthik Muthuraman, and Zachary Eichenberger. 2020. Identifying incorrect labels in the conll-2003 corpus . In <i>Proceedings of the 24th Conference on Computational Natural Language Learning, CoNLL 2020, Online, November 19-20, 2020</i> , pages 215–226. Association for Computational Linguistics.	1031
977			1032
978			1033
979			1034
980			1035
981			1036
982			1037
983	Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.	Simon Rogers, Derek H. Sleeman, and John Kinsella. 2013. Investigating the disagreement between clinicians' ratings of patients in icus . <i>IEEE J. Biomed. Health Informatics</i> , 17(4):843–852.	1038
984			1039
985			1040
986			1041
987		Tal Schuster, Adam Fisch, and Regina Barzilay. 2021. Get your vitamin C! robust fact verification with contrastive evidence . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 624–643, Online. Association for Computational Linguistics.	1042
988			1043
989			1044
990	An Thanh Nguyen, Byron C. Wallace, and Matthew Lease. 2015. Combining crowd and expert labels using decision theoretic active learning . In <i>AAAI Conference on Human Computation & Crowdsourcing</i> .		1045
991			1046
992			1047
993			1048
994			
995	Duc Tam Nguyen, Chaithanya Kumar Mummadi, Thi-Phuong-Nhung Ngo, Thi Hoai Phuong Nguyen, Laura Beggel, and Thomas Brox. 2019. Self: Learning to filter noisy labels with self-ensembling . <i>ArXiv</i> , abs/1910.01842.	Katharine Sherratt, Hugo Gruson, Rok Grah, Helen Johnson, Rene Niehus, Bastian Prasse, and et al. 2023. Predictive performance of multi-model ensemble forecasts of covid-19 across european nations . <i>eLife</i> , 12:e81916.	1049
996			1050
997			1051
998			1052
999			1053
1000	Curtis G. Northcutt, Anish Athalye, and Jonas W. Mueller. 2021. Pervasive label errors in test sets destabilize machine learning benchmarks . <i>ArXiv</i> , abs/2103.14749.	Rion Snow, Brendan T. O'Connor, Dan Jurafsky, and A. Ng. 2008. Cheap and fast – but is it good? evaluating non-expert annotations for natural language tasks . In <i>Conference on Empirical Methods in Natural Language Processing</i> .	1054
1001			1055
1002			1056
1003			1057
1004	Curtis G. Northcutt, Lu Jiang, and Isaac L. Chuang. 2019. Confident learning: Estimating uncertainty in dataset labels . <i>J. Artif. Intell. Res.</i> , 70:1373–1411.	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, and et al. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models . <i>Trans. Mach. Learn. Res.</i> , 2023.	1059
1005			1060
1006			1061
1007	OpenAI. 2023. GPT-4 technical report . <i>CoRR</i> , abs/2303.08774.		1062
1008			1063
1009	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe.		1064
1010			1065
1011			1066
1012			1067
1013			1068
1014			1069
			1070
			1071

1072	Julius Steen, Juri Opitz, Anette Frank, and Katja Markert. 2023. With a little push, NLI models can robustly and efficiently predict faithfulness . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 914–924, Toronto, Canada. Association for Computational Linguistics.	1128
1073		1129
1074		1130
1075		1131
1076		1132
1077		1133
1078		1134
1079	Aneeta Sylolypavan, Derek H. Sleeman, Honghan Wu, and Malcolm Sim. 2023. The impact of inconsistent human annotations on AI driven clinical decision making . <i>npj Digit. Medicine</i> , 6.	1135
1080		1136
1081		1137
1082		1138
1083	Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. 2016. Re-thinking the inception architecture for computer vision . In <i>2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016</i> , pages 2818–2826. IEEE Computer Society.	1139
1084		1140
1085		1141
1086		1142
1087		1143
1088		1144
1089		1145
1090	Derek Tam, Anisha Mascarenhas, Shiyue Zhang, Sarah Kwan, Mohit Bansal, and Colin Raffel. 2023. Evaluating the factual consistency of large language models through news summarization . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 5220–5255, Toronto, Canada. Association for Computational Linguistics.	1146
1091		1147
1092		1148
1093		1149
1094		1150
1095		1151
1096		1152
1097	James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and VERification . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 809–819, New Orleans, Louisiana. Association for Computational Linguistics.	1153
1098		1154
1099		1155
1100		1156
1101		1157
1102		1158
1103		1159
1104		1160
1105		1161
1106	Petter Törnberg. 2023. Chatgpt-4 outperforms experts and crowd workers in annotating political twitter messages with zero-shot learning . <i>ArXiv</i> , abs/2304.06588.	1162
1107		1163
1108		1164
1109		1165
1110	Alexandra Uma, Tommaso Fornaciari, Dirk Hovy, Silviu Paun, Barbara Plank, and Massimo Poesio. 2021. Learning from disagreement: A survey . <i>J. Artif. Intell. Res.</i> , 72:1385–1470.	1166
1111		1167
1112		1168
1113		1169
1114	Mor Ventura, Eyal Ben-David, Anna Korhonen, and Roi Reichart. 2023. Navigating cultural chasms: Exploring and unlocking the cultural POV of text-to-image models . <i>CoRR</i> , abs/2310.01929.	1170
1115		1171
1116		1172
1117		1173
1118	Veniamin Veselovsky, Manoel Horta Ribeiro, Philip Cuzzolino, Andrew Gordon, David Rothschild, and Robert West. 2023a. Prevalence and prevention of large language model use in crowd work . <i>CoRR</i> , abs/2310.15683.	1174
1119		1175
1120		1176
1121		1177
1122		1178
1123	Veniamin Veselovsky, Manoel Horta Ribeiro, and Robert West. 2023b. Artificial artificial intelligence: Crowd workers widely use large language models for text production tasks . <i>CoRR</i> , abs/2306.07899.	1179
1124		1180
1125		1181
1126		1182
1127		1183
		1184
		1185
		1186
	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. GLUE: A multi-task benchmark and analysis platform for natural language understanding . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.	1187
		1188
		1189
		1190
		1191
		1192
		1193
		1194
		1195
		1196
		1197
		1198
		1199
		1200
		1201
		1202
		1203
		1204
		1205
		1206
		1207
		1208
		1209
		1210
		1211
		1212
		1213
		1214
		1215
		1216
		1217
		1218
		1219
		1220
		1221
		1222
		1223
		1224
		1225
		1226
		1227
		1228
		1229
		1230
		1231
		1232
		1233
		1234
		1235
		1236
		1237
		1238
		1239
		1240
		1241
		1242
		1243
		1244
		1245
		1246
		1247
		1248
		1249
		1250
		1251
		1252
		1253
		1254
		1255
		1256
		1257
		1258
		1259
		1260
		1261
		1262
		1263
		1264
		1265
		1266
		1267
		1268
		1269
		1270
		1271
		1272
		1273
		1274
		1275
		1276
		1277
		1278
		1279
		1280
		1281
		1282
		1283
		1284
		1285
		1286
		1287
		1288
		1289
		1290
		1291
		1292
		1293
		1294
		1295
		1296
		1297
		1298
		1299
		1300

- 1187 [risk minimization](#). In *6th International Conference*
1188 *on Learning Representations, ICLR 2018, Vancouver,*
1189 *BC, Canada, April 30 - May 3, 2018, Conference*
1190 *Track Proceedings*. OpenReview.net.
- 1191 Ruoyu Zhang, Yanzeng Li, Yongliang Ma, Ming Zhou,
1192 and Lei Zou. 2023. [Llmeta: Making large language](#)
1193 [models as active annotators](#). *ArXiv*, abs/2310.19596.
- 1194 Yuan Zhang, Jason Baldridge, and Luheng He. 2019.
1195 [PAWS: Paraphrase adversaries from word scrambling](#).
1196 In *Proceedings of the 2019 Conference of the North*
1197 *American Chapter of the Association for Computa-*
1198 *tional Linguistics: Human Language Technologies,*
1199 *Volume 1 (Long and Short Papers)*, pages 1298–1308,
1200 Minneapolis, Minnesota. Association for Computa-
1201 tional Linguistics.
- 1202 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
1203 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
1204 Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,
1205 Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging](#)
1206 [llm-as-a-judge with mt-bench and chatbot arena](#). In
1207 *Advances in Neural Information Processing Systems*
1208 *36: Annual Conference on Neural Information Pro-*
1209 *cessing Systems 2023, NeurIPS 2023, New Orleans,*
1210 *LA, USA, December 10 - 16, 2023*.

Appendix

A Additional Experiments - SummEval	15
A.1 Data	15
A.2 Definitions	15
A.3 Experimental Setting	15
A.4 Experiments and Results	16
B The Power of Ensemble	16
C Comparing Annotation Approaches	17
C.1 Annotation Quality	17
C.2 Consistency	18
C.3 Cost and Scalability	19
D Annotation	20
D.1 Crowd-source	20
D.2 Experts	22
D.3 LLMs	22
E Data	23
F Misabeled Data Implications	25
F.1 Fine-tuning	25
F.2 Model Evaluation	25
G Statistical Analysis	25
G.1 Clopper-Pearson	25
G.2 Bootstrap sampling	26
H Label Errors	26

A Additional Experiments - SummEval

In addition to the datasets from the TRUE benchmark, we replicate our experiments on another dataset with a different objective and a different labeling scheme, to strengthen our results and conclusions.

A.1 Data

SummEval (Fabbri et al., 2021) is an extensive and commonly used summarization benchmark, evaluating the quality of multiple model-generated summarization outputs compared to a source CNN/DailyMail sources on four dimensions: coherence, relevance, consistency, and fluency. Each summarization is labeled on each dimension with five crowd-workers and three experts,

enabling us to replicate some of the experiments without additional crowd-worker or expert annotation costs. The labeling schema is ordinal on a scale of 1 to 5 (higher is better). Note that this dataset does not have a singular gold-standard label per summarization, but rather a collection of annotations from experts and crowd-workers. Therefore, we will not claim to find label errors in this benchmark, but rather showcase our methodology as if the crowd-sourced annotations are the original labels for the dataset, and we have access to experts' annotations for gold-standard reference, to determine if the LLM was correct when flagging examples.

A.2 Definitions

To apply our methods for error detection via LLMs ensemble, we first define the following:

Labels We aggregate crowd-sourced annotations by their median, to construct a single original label on a scale of 1 to 5. Similarly, we take the median of the experts' annotations to be a single gold-standard label.

A disagreement We say that the LLM annotation *disagrees* with the original label if there is a difference of more than 1 between the scores. The idea is that we can confidently say the LLM and the original label "disagree", as if the difference is 1 or less, this is a weak disagreement we will probably not flag for.

A.3 Experimental Setting

Similar to the description in subsection 4.2, we utilize two LLMs— GPT-4o (gpt-4o-2024-11-20) and Gemini 1.5 Flash (gemini-1.5-flash-002). We constructed four prompts, differing by phrasing and compatible with the four prompt template structures used for the TRUE benchmark experiments. The answer to each query was a JSON format with 'Relevance', 'Coherence', 'Consistency', and 'Fluency' as its keys. The scores are integers on a scale of 1 to 5, as are the ratings in the SummEval dataset. We extract the probability of each score possible through the log-probs for each score token. Finally, we average all models' probabilities, to obtain an ensemble of LLMs, with p being the distribution over the five possible scores.

1301	A.4 Experiments and Results	
1302	A.4.1 Can LLMs Detect Label Errors?	
1303	We replicate the experiment described in subsection 5.2 with the appropriate adjustment for the SummEval dataset, based on the definitions above.	
1304	The result is shown in Figure 2 (bottom). The plot presents the subset of examples where there was a disagreement between the crowd-sourced annotation and the LLMs' annotation. Each bin represents the confidence of the LLMs in their predicted label.	
1305	As there are five ordinal categories, even if there was a disagreement between two annotations, they both might be "wrong", where the expert's answer is a third option. Therefore, to show clearer results, we do not resolve by experts "who is correct", but rather "who is more correct?". For completeness, we also provide the "both equally correct" option, for the case the expert's label is exactly in the middle, and none is "more correct" than the other. The bins are relatively balanced in terms of the amount of examples per bin. Note that in contrast to the TRUE binary labeling scheme, where confidence 0.5 is the minimal threshold for an answer, here we start from 0.2.	
1306		
1307		
1308		
1309		
1310		
1311		
1312		
1313		
1314		
1315		
1316		
1317		
1318		
1319		
1320		
1321		
1322		
1323		
1324		
1325		
1326		
1327		
1328		
1329		
1330		
1331		
1332		
1333		
1334		
1335		
1336		
1337		
1338		
1339	A.4.2 The Power of Ensemble	
1340	We analyze the importance of utilizing more than a single model and a single prompt on two dimensions - performance compared to the gold labels (the quality of the annotations we utilize), and error detection (the ability to identify errors more accurately). For performance evaluation on the ordinal labels, we report Pearson correlation; for error detection evaluation, we report the F1-score based on binary error/not-error classification. See results in Figure 3 and discussion in Appendix B .	
1341		
1342		
1343		
1344		
1345		
1346		
1347		
1348		
1349		
	A.4.3 Annotation Approaches Comparison	1350
	In Appendix C , we thoroughly discuss the comparison between the different annotation approaches. For SummEval, experts and crowd-sourced annotations are provided. Together with our LLM-ensemble annotations (as described in subsection A.3), we analyze and compare the annotation approaches in terms of quality (see Figure 6 (bottom)) and consistency (see Table 5). To account for ordinal labels, we measure IAA via Krippendorff's α (Krippendorff, 1970).	1351 1352 1353 1354 1355 1356 1357 1358 1359 1360
	B The Power of Ensemble	1361
	As mentioned in subsection 4.2 , we treat the LLM annotations as an ensemble of 2 models combined with 4 different prompts, in order to ensure greater stability in the results. Where one LLM may succeed, the other may fail, and averaging all their probabilities enables us to have more confidence in the final answer. In this subsection, we further analyzed the performance of LLMs by varying the size of the LLM ensemble, examining how this affects the model performance. We evaluate two aspects of model performance. First, we assess how closely the ensemble's annotations match the gold labels— essentially, how much we can trust the LLM annotations. We measure this aspect of label quality using the ROC AUC compared to the gold labels. The second aspect is the ensemble's ability to detect label errors. For this, we compute the F1-score by averaging the recall of errors and the precision of correctly identifying a candidate as a true error.	1362 1363 1364 1365 1366 1367 1368 1369 1370 1371 1372 1373 1374 1375 1376 1377 1378 1379 1380 1381
	Results are shown in Figure 3 (top). For both aspects, we see a clear trend. As we increase the number of models in the ensemble, the performance increases. In terms of ROC AUC w.r.t the gold labels (left plot), this suggests better annotation quality, while the right plot, a higher F1 score indicates a stronger error detector, either by recalling more errors or improving precision, or through a balance of both. Additionally, for both measures, the variance decreases as the ensemble size grows, which indicates more stable and consistent annotations and error detections. Similarly, Figure 3 (bottom) shows the power of LLM ensemble on the same aspects on the SummEval datasets, aggregated over four summarization dimensions (see experiment details on Appendix A.4.2). Trends of diminishing variance and increased performance and error detection are observed here as well.	1382 1383 1384 1385 1386 1387 1388 1389 1390 1391 1392 1393 1394 1395 1396 1397 1398 1399

Although not yet discussed in the context of error detection with LLMs, these results align with previous work showing the power of ensemble (Dietterich, 2007). These observations justify our choice to use an ensemble of models rather than a single one.

C Comparing Annotation Approaches

Our paper discusses three annotation approaches, each with its own benefits and drawbacks. These approaches differ in how they manage the trade-offs between label quality, scalability, and cost. In the following section, we discuss and compare their characteristics. A summary of this comparison is given in Figure 4.

C.1 Annotation Quality

When annotating or validating a dataset, one of our main concerns is the quality of the labels, or in other words, establishing a reliable gold standard. However, each annotation approach produces different labels. To estimate the quality of these approaches, we measure the agreement between different annotations using the weighted F1-score (which accounts for both classes). Note that this metric is not symmetric, meaning that treating one annotation as the *true* label and the other as the *prediction*, or vice versa, can result in different scores.

Figure 6 (top) presents the F1-score between each pair of annotation approaches. As the figure shows, LLMs have disagreements with the *original* labels (0.72). Yet, as discussed in subsection 5.1, the original labels themselves contain mistakes, so this disagreement does not necessarily indicate poor performance of the LLMs. When considering the *Gold* as the true label, LLM performance increases to 0.83. This suggests that LLMs, despite their discrepancies with the original labels, perform closer to the truth than initially reported. The *Gold* label, obtained by experts, has high agreement with both the *Original* and *LLM* labels. On the other hand, the *MTurk-Majority* approach performs poorly, with near-random F1-scores compared to both the original and gold labels, and even when compared to its stricter variant, *MTurk-Strict*. The results indicate that basic crowd-sourcing, without additional training to enhance crowd-workers into specialized sub-experts, performs significantly worse compared to other approaches, including LLM-based methods. On the SummEval dataset

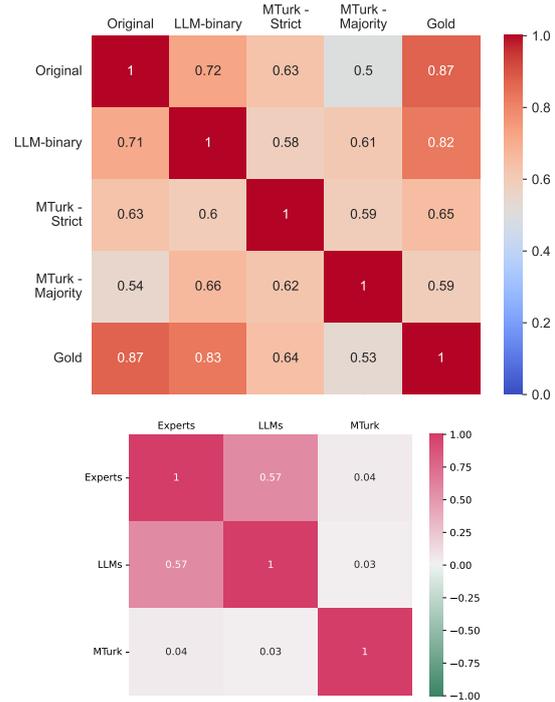


Figure 6: Comparison between all annotation methods: **(Top)** on the TRUE benchmark, measured by the weighted-F1-score. Rows represent the *"true"* label and columns represent the *"prediction"*. For instance, the score of *LLMs* compared to the *Original* label is 0.72. **(Bottom)** Comparison on the SummEval benchmark, measured by Pearson correlation (results are averaged over all dimensions).

(bottom of Figure 6 bottom) we observe similar results, where the LLMs are more correlated with the Experts rather than the crowd-workers, which in turn have almost-no-correlation with LLMs or experts' annotations— this implies poor quality of the annotations obtained from crowd-source.

Crowd-sourcing For crowd-sourcing, the reported F1-score does not provide the complete picture. When we focus on individual annotators, we see that those who annotate more examples generally deliver higher-quality annotations, achieving greater accuracy when compared to both the original and gold labels (see Figure 7). This phenomenon can be explained by two hypotheses: (1) a learning process— as the annotators see more examples, they improve at the task, or (2) users who dedicate time to annotating multiple examples are likely those who either read the guidelines carefully and strive to perform the task to the best of their ability, or are naturally proficient at the task and therefore continue annotating. Even though annotators who label more instances tend to provide higher-quality

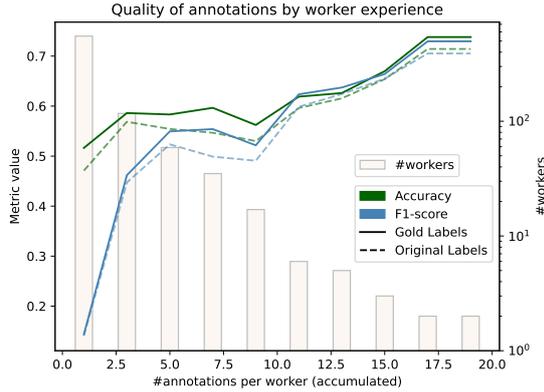


Figure 7: (x-axis) at list x annotations per annotator. (Right y-axis) The number of annotators with at least x annotations (bins). (Left y-axis) the average F1-score or accuracy for all user annotations with at least x annotations.

1471 annotations, they are less common—most annotators
 1472 tend to stop after only a few examples. This dis-
 1473 tribution of annotators results in overall insufficient
 1474 annotation quality. Pre-qualification tests are of-
 1475 ten used to shift this distribution from the "average
 1476 worker" towards more experienced or dedicated
 1477 annotators; however, this requires a significantly
 1478 larger budget and greater micro-management in-
 1479 volvement from the researcher.

1480 C.2 Consistency

1481 Usually, when annotating a dataset, more than
 1482 one annotator is involved. This applies to crowd-
 1483 workers, experts, and even LLMs—in this study,
 1484 we use an ensemble of different LLMs and prompts.
 1485 The use of multiple annotators, similar to an en-
 1486 semble, is meant to overcome the variance between
 1487 individuals, which can arise from the subjective
 1488 nature of NLP tasks, different interpretations of
 1489 instructions, lack of experience, task difficulty, and
 1490 cognitive bias (Uma et al., 2021).

1491 As such, a common practice in the NLP commu-
 1492 nity is to report Inter Annotator Agreement (IAA)—
 1493 a set of statistical measures used to evaluate the
 1494 agreement between individuals. Typically, IAA
 1495 can be viewed as an adjustment of the proportion
 1496 of pairwise agreements, where 0.0 indicates ran-
 1497 dom agreement. We focus on Fleiss’s κ (Fleiss,
 1498 1971), as it accounts for label imbalance and mul-
 1499 tiple (> 2) annotators. High IAA, or low variance
 1500 between independent annotators, is considered an
 1501 indicator of high-quality annotation. In Table Ta-
 1502 ble 4, we report the agreement between annotators
 1503 across different approaches. For LLMs, we report

two variants: (1) same model, different prompts;
 and (2) different models, where each model’s result
 is the aggregation across prompts. For reference,
 we also include the IAA from the original annota-
 tions, as reported in the original papers: *MNBM*
 reported an average Fleiss’s κ of 0.696 for the hal-
 lucination annotation task; *BEGIN* reported Krip-
 pendorff’s α (a generalization of Fleiss’s κ) of 0.7;
VitaminC reported Fleiss’s κ of 0.7065 on a sample
 of 2,000 examples; and *PAWS* reported a 94.7%
 agreement between a single annotator’s label and
 the majority vote on the Wikipedia subset used in
 TRUE.

Experts While it’s true that reconciliation natu-
 rally leads to increased agreement, the significant
 improvement in IAA we observed highlights its im-
 portance. Though this phase is less common in
 practice, it is crucial not only for increasing agree-
 ment but also for improving the overall quality of
 annotations and ensuring more reliable outcomes.
 Interestingly, label changes in this phase were not
 symmetric, as most changes (69.3%) were in the
 direction of *consistent* \rightarrow *inconsistent*, where one
 annotator found an inconsistency that the other did
 not (see all change details in Figure 11). It is impor-
 tant to note that the κ obtained by the experts (both
 before and after reconciliation) was calculated on
 a more challenging subset, where the original label
 differed from the LLM prediction, and should be
 interpreted with this context in mind. This is re-
 flected in the decrease in κ observed for all other
 annotator groups on this subset.

LLMs GPT-4 and PaLM2, the better-performing
 LLMs on this task, show high IAA, with $\kappa = 0.706$
 and $\kappa = 0.75$, respectively, which is similar to the
 experts’ reported κ . This suggests a comparable
 level of variance and quality in annotation, pro-
 viding further empirical evidence for considering
 LLMs as annotators. This property adds to previ-
 ous studies showing LLMs’ quality as surrogates
 for human preferences (Zheng et al., 2023) or eval-
 uations (Chiang and Lee, 2023).

Crowd-Sourcing. Crowd workers showed near-
 random agreement, indicating relatively poor-

*Multiple MTurk workers have participated in annotation, yet exactly 3 annotations per example were obtained. Annotators independence assumption was made to calculate Fleiss’s κ as with 3 annotators.

†These MTurk annotators were chosen with stricter pre-qualification criteria than those in the TRUE dataset and do not correspond to the MTurk line in the TRUE table.

Annotator group	Fleiss's κ	%agreement	#examples	Fleiss's κ (disagree. subset)	#annotators
Experts			222		2
Before reconciliation	0.486	75.7		0.486	
After reconciliation	0.851	93.2		0.851	
MTurk	0.074	60.5	400	-0.004	3*
LLM (different prompts)			640		4
GPT-4	0.706	85.3		0.571	
PaLM2	0.750	87.7		0.696	
LLaMA3	0.219	71.7		0.078	
Mistral	0.459	73.2		0.314	
LLMs (different models)	0.521	77.5	640	0.389	4

Table 4: Inter-Annotator Agreement in different annotator groups. %agreement is the proportion of pairwise annotator comparison. Fleiss's κ (disagree. subset) refers to the κ over the subset of disagreement between LLM and the original label.

Annotator group	Krippendorff's α	%agreement	#annotators
Experts	0.584	60.4	3
MTurk[†]	0.496	65.6	5
LLM (different prompts)			4
GPT-4o	0.760	63.6	
Gemini 1.5 Flash	0.733	79.7	
LLMs (different models)	0.576	62.9	2

Table 5: Inter-Annotator Agreement in different annotator groups on the SummEval benchmark. %agreement is the proportion of pairwise annotator comparisons.

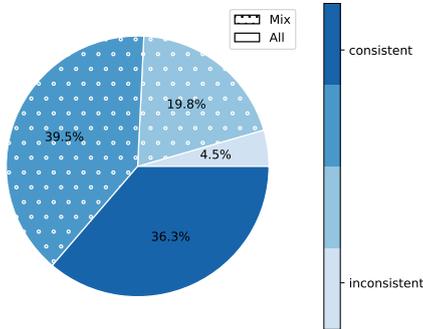


Figure 8: Distribution of crowd-source annotators. Each example was annotated by 3 workers. Plain segments are unanimous annotation, while dotted segments indicate examples where some annotators labeled as *inconsistent*, and other as *consistent*. For example, 19.8% of the examples had two *inconsistent* annotation, and one *consistent* annotation.

quality annotations. Figure 8 describes the distribution of annotations by MTurk workers. Only 40.8% of the examples were labeled unanimously, whereas the rest included annotations from both classes. In addition, if aggregating by majority vote, we get that 75.8% of the examples are labeled as *consistent*, which is far from the original distribu-

tion of classes. As mentioned before, even experts may miss a small inconsistency nuance, and finding it requires attention. Even from the subset of examples unanimously labeled as *consistent*, 37.9% have a label of *inconsistent* in both original and gold labels, which points to a lack of attention and thoroughness.

SummEval. Table 5 shows the IAA analysis on the SummEval benchmark. We report Krippendorff's α (Krippendorff, 1970), a generalization of κ to account for ordinal labeling. LLMs exhibit high IAA (compared to experts' IAA) of $\alpha = 0.57$ and 62.9% agreement between models, with high consistency across prompts for the same model. Crowd-workers obtain decent results (maybe due to stricter pre-qualification criteria of 10,000 approved HITs), yet they still fall short compared to experts or LLMs.

C.3 Cost and Scalability

In MTurk platform, a total of $400 \times 3 = 1200$ annotations cost 572\$, including 2 small pilot experiments. All annotations were prepared within a few hours. However, it demanded an additional and significant time for review, after which rejected exam-

1579 ples returned to the pool. This annotation-review
1580 cycle was conducted for ~ 5 iterations. Infer-
1581 ence via OpenAI’s API on GPT-4 cost $\sim 4.5\%$ per
1582 prompt. Inference via VertexAI’s API on PaLM2
1583 cost $\sim 0.15\%$ per prompt. Both took ~ 8 minutes
1584 per prompt. Inference on Mistral and Llama3
1585 was via the HuggingFace API, and its cost is esti-
1586 mated by the cost of using a suitable Virtual Ma-
1587 chine (VM) on Google Cloud Platform (GCP) for
1588 the time of inference (1 minute per model)- $\sim 0.1\%$
1589 per prompt.

1590 LLM-based annotation is significantly cheaper
1591 and faster than crowd-sourcing platforms like
1592 MTurk, especially when considering the additional
1593 time required for human review cycles. It is esti-
1594 mated to be 100 to 1,000 times more cost-effective
1595 than using human annotators, including experts.
1596 This scalability and speed make LLMs a highly ef-
1597 ficient alternative for large-scale annotation tasks.

1598 **D Annotation**

1599 **D.1 Crowd-source**

1600 Each example was annotated by three annotators,
1601 who in addition to the binary label were requested
1602 to provide their confidence in their answer, and also
1603 write a short explanation for why they chose this la-
1604 bel. Pre-qualifications included 50+ approved HITs
1605 and 97%+ approval rate, which are at standard scale
1606 for the MTurk platform (Kazai et al., 2013; Hauser
1607 et al., 2021; Chmielewski and Kucker, 2019). Also,
1608 locations were limited to [USA, UK, Australia],
1609 which are all English-speaker countries. We dis-
1610 abled the possibility of right-click and `Ctrl+C` in
1611 the platform (as suggested by (Veselovsky et al.,
1612 2023a)), to prevent (as much as possible) the case
1613 where generative-AI (e.g., ChatGPT) will be ap-
1614 plied to solve the task instead of humans solv-
1615 ing it themselves (as shown by (Veselovsky et al.,
1616 2023b)). The maximum time allowed per HIT was
1617 6 minutes, while the actual average execution time
1618 was 2:20 minutes for all assignments, and 3 min-
1619 utes for approved assignments. The guidelines pro-
1620 vided to annotators and the annotation platform
1621 layout are presented in Figure 9.

1622 Each annotation was manually reviewed and was
1623 rejected if the answers were not in line with the in-
1624 structions, or if it was obvious that the task was not
1625 done honestly. Overall, this task suffered from a
1626 high rejection rate of 49.2% (1163 rejected, 1200
1627 approved). The main rejection reasons were: lack
1628 of meaningful explanation, obvious copy-paste an-

notations across different examples, explanations
1629 contradicting the label annotation, and cases where
1630 the explanation was a copy-paste of either the
1631 grounding or the statement. 1632

Factual Consistency Evaluation - Instructions ✕

Thank you for participating in our research on factual consistency in texts.

Each example consists of two texts:

1. **Grounding** - A factual text.
2. **Statement** - A text to be evaluated.

Task:

Your task is to determine if the Statement is factually consistent with the Grounding.

Definition of Factual Consistency:

- **Factual Consistency:** The Statement accurately reflects and aligns with all the facts presented in the Grounding. The Statement does not introduce any errors, new entities, or unsupported information and is in full agreement with the Grounding.
- **Factual Inconsistency:** The Statement contains any inaccuracies, contradictions, or information that cannot be supported by the Grounding or derived from it.

Answer Format:

Your answer should be binary: either **Factually Consistent** or **Factually Inconsistent** (choose the appropriate answer in the "Your Answer" section).

Additional Information Required:

- Confidence Level: Indicate your confidence in your answer on a scale of 1 to 5 ("Your Confidence").
- Explanation: Provide a brief explanation for your answer ("Short Explanation" text box).

We appreciate your attention to detail and accuracy in this evaluation process. Thank you for your valuable contribution.

Grounding:

At the same time , Pope Francis Tong asked Bishop of Hong Kong to stay for three years .

Statement:

At the same time , Pope Francis asked Tong to remain Bishop of Hong Kong for three more years .

Your task is to determine if the Statement is factually consistent with the Grounding.

Your Answer:

Factually Inconsistent
 Factually Consistent

Your Confidence:

Indicate your confidence in your answer on a scale of 1 to 5.
(Note: 0 is not part of the scale)

Short Explanation:

Provide a brief but meaningful explanation (at least one sentence) for why you classified the statement as factually consistent or inconsistent

Figure 9: Platform for crowd-sourcing annotation in Amazon Mechanical Turk (MTurk). **(Top)** Guidelines for the task and definitions. **(Bottom)** Annotation layout for a single instance.

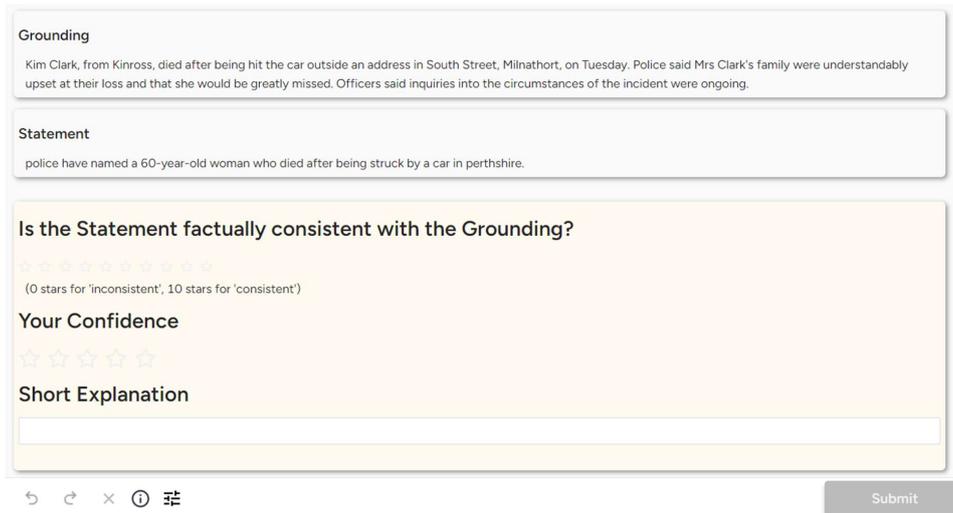


Figure 10: Annotation platform on Label-Studio for experts

D.2 Experts

Experts annotation was using the platform of Label Studio.¹ Layout design is presented in Figure 10. Examples were presented in random order, and neither the LLM prediction nor the original label were presented during the annotation. In the first stage, each example was annotated independently by both experts. Afterward, the human experts began in a second phase of a reconciliation, where a discussion was made over examples they disagreed over. This reconciliation phase ended up with a much higher agreement and higher-quality labels.

In the reconciliation phase, we observed that most changes (69.3%) were from label 1 to label 0, indicating that contradictions might be hard to find, and not all annotators catch them at first. For the full distribution of label change in the reconciliation phase, see Figure 11.

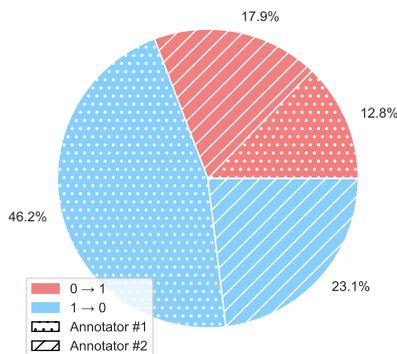


Figure 11: How experts' annotations have changed after the reconciliation phase. Most changes occur from 1 (*consistent*) to 0 (*inconsistent*).

¹<https://labelstud.io/>

D.3 LLMs

To annotate a total of $160 \times 4 = 640$ examples from four different datasets, we used four LLMs: GPT-4 (gpt-4-1106-preview) (OpenAI, 2023), PaLM2 (text-bison@002) (Anil et al., 2023), Mistral (7B)² (Jiang et al., 2023) and Llama 3 (8B)³ (Dubey et al., 2024).

Each model was run with four different prompts (see full prompts in Figure 12). We used a variety of terminology, as this task appears to have different framings in different studies. For example, the premise-hypothesis terminology from classic NLI (MacCartney and Manning, 2009), or document-statement used in (Tam et al., 2023).

For API models (GPT-4, PaLM2), we set `temperature=0.0` and extracted the logit of the generated token (functionality provided by both APIs), if the generated token was either '0' or '1' as expected. This logit was then transformed into a probability $p_t = P(y = t|x)$ via exponent corresponding the generated token t , and $1 - p_t$ for the other label. To address the case where the first generated token was an unrelated token such as ' ', '\n', we set `max_tokens=2` and took the first appearance of either '0' or '1'. For all models, prompts and examples, '0' or '1' were in the first two generated tokens. Rest of parameters were set according to their default values.

For models available through the HuggingFace API (e.g., Mistral, Llama 3), we can load the model

²<https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.2>

³<https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct>

parameters and make inference locally. In that case, we get access to logits for all tokens, instead of just for the generated ones. Therefore, we applied a similar procedure, where we seek for the first appearance of either '0' or '1' to be the most probable token to be generated, and then directly extracted the logits of the '0' and '1' tokens. These logits were transformed into probabilities ($P(y = 0|x)$, $P(y = 1|x)$) via a softmax function.

E Data

For our main experiments, we used the TRUE benchmark for factual consistency. Specifically, we focus on four TRUE datasets, one from each task (summarization, dialogue, fact verification, paraphrasing):

MNBM (Maynez et al., 2020): Summarization.

This dataset provides annotations for hallucinations in generated summaries from the XSum dataset (Narayan et al., 2018). *Grounding* refers to the source document that the summary is based on, while *Generated Text* consists of model-generated summaries, which may include hallucinated information not present in the source. Three human annotators, trained for the task through two pilot studies, annotated the dataset for the existence of hallucinations. In TRUE, the binary annotations were determined by majority vote.

BEGIN (Dziri et al., 2022): Dialogue. This dataset evaluates groundedness in knowledge-grounded dialogue systems, where responses are expected to align with an external *Grounding* source, typically a span from Wikipedia. *Generated Text* refers to model-generated dialogue responses that were fine-tuned on datasets like Wizard of Wikipedia (Dinan et al., 2019). Data was annotated into entailment/neutral/contradiction labels, by three human annotators, trained for the task through two pilot studies, aggregated by majority vote. In TRUE, binary annotations were then determined by the entailment/not-entailment partition.

VitaminC (Schuster et al., 2021): Fact Verification. This dataset is based on factual revisions of Wikipedia. The evidence, or *Grounding*, consists of Wikipedia sentences, either before or after these revisions. Most human involvement came from creating *Generated Text* rather than the annotation process, with annotators writing claim/evidence pairs derived from Wikipedia revisions, inherently

generating labeled data for fact verification. Synthetic examples from the FEVER dataset (Thorne et al., 2018) were also included. Additionally, three annotators reviewed 2,000 examples, presumably to ensure data quality.

PAWS (Zhang et al., 2019): Paraphrasing.

This dataset consists of paraphrase and non-paraphrase pairs. *Grounding* refers to source sentences drawn from Quora and Wikipedia, while *Generated Text* was automatically generated through controlled word swapping and back-translation. Five human annotators annotated the dataset with binary labels w.r.t paraphrasing correctness. The dataset includes both high- and low-agreement annotations.

1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743

prompt1

Here are two texts:

TEXT 1. <..PREMISE..>.

TEXT 2. <..HYPOTHESIS..>.

Is TEXT 2 contradictory or is it factually inconsistent with TEXT 1? If yes answer 0.

Is TEXT 2 entailed or is it factually consistent with TEXT 1? If yes answer 1.

Refer only to the two texts above, and not any other external knowledge or context.

Answer only 0 or 1

Answer only with one token: 0 or 1

Answer:

prompt2

DOCUMENT: <..PREMISE..>.

QUESTION: Is the following STATEMENT factually consistent with the above document?

STATEMENT: <..HYPOTHESIS..>.

ANSWER FORMAT: 0 for No, 1 for Yes

Answer only with one token: 0 or 1

Answer:

prompt3

You are given the two following texts:

TEXT 1. <..PREMISE..>.

TEXT 2. <..HYPOTHESIS..>.

TEXT 1 is a fact. TEXT 2 is a statement. Is TEXT 2 factually consistent with TEXT 1?

Answer 0 for No, 1 for Yes.

Answer only with one token: 0 or 1

Answer:

prompt4

Given the following texts:

<PREMISE> : <..PREMISE..>.

<HYPOTHESIS> : <..HYPOTHESIS..>.

Please assess the factual consistency of <HYPOTHESIS> with respect to <PREMISE>.

If the content of <HYPOTHESIS> aligns with the information provided in <PREMISE>, assign a label of 1.

If there are factual inconsistencies between <HYPOTHESIS> and <PREMISE>, assign a label of 0.

Target Format: either 0 (for Factual Inconsistency) or 1 (for Factual Consistency).

Answer only with one token: 0 or 1

Answer:

Figure 12: Four different prompt input templates to LLMs for obtaining binary labels

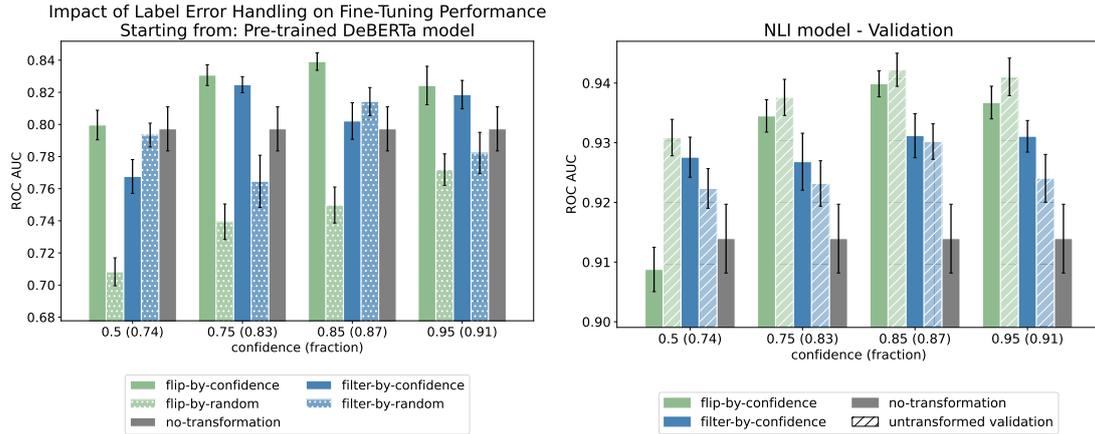


Figure 13: Similar experiments to the one in Figure 5, with small alterations. **(Left)** Starting from a different base model - pre-trained DeBERTa-v3-base. **(Right)** Dashed columns present results for when flipping or filtering methods were applied only on the training set, but not the validation.

F Mislabeled Data Implications

F.1 Fine-tuning

Hardware. For the finetuning of DeBERTa models, both the base pre-trained model, and the NLI model which is in the same size, in subsection 7.1, we used 2 Quadro RTX6000 (24GB) GPUs.

Implementation. We finetuned starting from two base models: DeBERTa-v3⁴, and a fine-tuned version of it on classic NLI datasets⁵. We used HuggingFace trainer with early stopping of 4 epochs. The finetuning procedure includes splitting the training set into train and validation sets (where validation size is 25% and train 75%), fine-tuning on the train set, and choosing the best checkpoint based on the validation ROC AUC. We ran all experiments on five different seeds, affecting also the train-validation split and the random set chosen for ablation. We fine-tuned all variants with the same hyperparameters, determined by the best performing on the no-manipulation baseline. This includes 30 epochs at most, batch size of 16, learning rate of $5e-5$ and weight-decay of 0.03. The rest were set as the trainer and model default.

Additional Experiments. The left plot in Figure 13 presents the same experiment discussed in subsection 7.1, but starting from the pre-trained DeBERTa-v3-base. Same trends applies here, where our LLM-confidence-based manipulations of either flipping or filtering flagged examples outperforms the baselines.

⁴microsoft/deberta-v3-base

⁵MoritzLaurer/DeBERTa-v3-base-mnli-fever-anli

The right plot in Figure 13 compares the performance of these methods (starting from the NLI model) when applied to both the training and validation sets (solid bars) or only the training set (dashed bars). The results are consistent, with no statistically significant differences between the two settings. Importantly, all variations outperform the baseline, underscoring the critical role of a well-curated training set in enhancing the model’s ability to generalize effectively.

F.2 Model Evaluation

In subsection 7.2 we evaluated the following models: GPT-4, PaLM2 (text-bison@002), Mistral-v0.2 (7B), and Llama3 (8B), which are covered in subsection 4.2; DeBERTa-v3 and NLI-model, which is a fine-tuned version of it on NLI datasets, as discussed in subsection 7.1; and GPT-4o, GPT-4o-mini, Mistral-v0.3,⁶ which share the same implementation as GPT-4 or Mistral-v0.2.

G Statistical Analysis

G.1 Clopper-Pearson

As mentioned in subsection 5.1, we employed the Clopper-Pearson exact method (Clopper and Pearson, 1934) to construct a 95% confidence interval for the binomial proportion, adjusted by a finite population correction (FPC). As we only have a subset of examples we re-annotated by LLMs or experts, we can not precisely determine what is the error rate in the full dataset, but only construct a confidence interval based on the re-annotated

⁶<https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3>

1804 subset. The Clopper-Pearson method provides an
1805 exact confidence interval for a binomial proportion,
1806 which means it gives a reliable estimate even with
1807 small sample sizes. By applying FPC, we adjust
1808 the interval because our sample is drawn from a
1809 limited population. This adjustment helps refine
1810 the estimate by taking into account the size of the
1811 overall dataset compared to the sample.

1812 **G.2 Bootstrap sampling**

1813 In [subsection 5.1](#), we use bootstrap sampling to
1814 provide confidence intervals for each bin. While
1815 not necessarily the first to introduce it, ([Xia et al.,
1816 2012](#)) explored bootstrap confidence intervals on
1817 ROC AUC. Unlike the method in [Appendix G.1](#),
1818 we do not make claims about the entire dataset,
1819 but rather focus on the re-annotated subset we pos-
1820 sess. To achieve this, we perform 100 bootstrap
1821 samples from the empirical distribution of each bin,
1822 sampling with replacement. We then measure the
1823 agreement between the experts' resolutions and the
1824 LLM annotations, compared to its agreement with
1825 the original label.

1826 **H Label Errors**

1827 [Table 6](#) demonstrates one example per dataset, in
1828 which the original label is, in fact, an error, the
1829 LLM prediction marked it as a candidate, and the
1830 expert annotators determined the correct gold label.
1831
1832

<p>Dataset: VITC</p> <p>Grounding: The British Government and NHS have set up a Coronavirus isolation facility at Arrowe Park Hospital in The Wirral for British People coming back on a special flight from Wuhan. Evacuation of foreign diplomats and citizens from Wuhan. Due to the effective lockdown of public transport in Wuhan and Hubei province , several countries have started to evacuate their citizens and/or diplomatic staff from the area , primarily through chartered flights of the home nation that have been provided clearance by Chinese authorities.</p> <p>Generated Text: There is a Coronavirus isolation facility at Arrowe Park Hospital that was set up by the NHS and the British Government</p> <p>Original Label: 0 LLM <i>p</i>: 0.99 Gold Label: 1</p> <p>Explanation: Rephrasing of the first sentence, without any contradiction.</p>
<p>Dataset: BEGIN</p> <p>Grounding: Hillary Clinton, the nominee of the Democratic Party for president of the United States in 2016, has taken positions on political issues while serving as First Lady of Arkansas (1979–81; 1983–92), First Lady of the United States (1993–2001);</p> <p>Generated Text: She is the nominee in 2016.</p> <p>Original Label: 0 LLM <i>p</i>: 0.98 Gold Label: 1</p> <p>Explanation: She (Hillary Clinton) is indeed the nominee in 2016 as specifically stated in the grounding.</p>
<p>Dataset: PAWS</p> <p>Grounding: David was born in Coventry on 21 September 1933 , with his twin Charles and Jessamine Robbins , the eighth and ninth children of twelve by Robbins.</p> <p>Generated Text: David was born on September 21 , 1933 in Coventry with his twin father Charles and Jessamine Robbins , the eighth and ninth child of twelve of Robbins</p> <p>Original Label: 1 LLM <i>p</i>: 0.04 Gold Label: 0</p> <p>Explanation: The generated text incorrectly states "twin father" instead of "twin" which is not the same, and does not even make much sense in English.</p>
<p>Dataset: MNBM</p> <p>Grounding: The John Deere tractor was pulled over by officers in the village of Ripley and had two other males on board. The vehicle had been seen in nearby Harrogate at about 05:00 GMT with no headlights on. Police said the driver had no licence, was not insured and did not have permission from the tractor’s owner. The vehicle was seized, with the three due to be interviewed by officers. Posting on Twitter, Insp Chris Galley said: "A strange end to a night shift. 15-year-old lad driving a tractor as a taxi for his drunk mates."</p> <p>Generated Text: a 15-year-old boy has been stopped by police after being seen driving a taxi on a night taxi.</p> <p>Original Label: 1 LLM <i>p</i>: 0.19 Gold Label: 0</p> <p>Explanation: The generated text claims that the 15-year-old boy was "driving a taxi on a night taxi", contradicting the grounding in which it was claimed that the boy was driving a tractor as a taxi</p>

Table 6: Annotation errors in the original datasets, discovered by LLMs and corrected by experts.