# FACTPICO: Factuality Evaluation for Plain Language Summarization of Medical Evidence

Sebastian Joseph<sup>1\*</sup> Lily Chen<sup>2\*</sup> Jan Trienes<sup>3,4</sup> Hannah Louisa Göke<sup>3,4</sup> Monika Coers<sup>3,4</sup> Wei Xu<sup>5</sup> Byron C. Wallace<sup>6</sup> Junyi Jessy Li<sup>1</sup>

<sup>1</sup>The University of Texas at Austin, <sup>2</sup>Massachusetts Institute of Technology <sup>3</sup>University of Duisburg-Essen, <sup>4</sup>Institute for AI in Medicine, University Hospital Essen <sup>5</sup>Georgia Institute of Technology, <sup>6</sup>Northeastern University

#### **Abstract**

Plain language summarization with LLMs can be useful for improving textual accessibility of technical content. But how factual are these summaries in a high-stakes domain like medicine? This paper presents FACTPICO, a factuality benchmark for plain language summarization of medical texts describing randomized controlled trials (RCTs), which are the basis of evidence-based medicine and can directly inform patient treatment. FACTPICO consists of 345 plain language summaries of RCT abstracts generated from three LLMs (i.e., GPT-4, Llama-2, and Alpaca), with fine-grained evaluation and natural language rationales from experts. We assess the factuality of critical elements of RCTs in those summaries: Populations, Interventions, Comparators, Outcomes (PICO), as well as the reported findings concerning these. We also evaluate the correctness of the extra information (e.g., explanations) added by LLMs. Using FACTPICO, we benchmark a range of existing factuality metrics, including the newly devised ones based on LLMs. We find that plain language summarization of medical evidence is still challenging, especially when balancing between simplicity and factuality, and that existing metrics correlate poorly with expert judgments on the instance level.

FactPICO and our code is available at https://github.com/lilywchen/FactPICO.

#### 1 Introduction

New findings in medicine observed in randomized controlled trials (RCTs) are published in journal articles which describe their design and outcomes. These RCTs "measure the effectiveness of a new intervention or treatment" (Hariton and Locascio, 2018) and are the important basis of evidence-based medicine (Sackett et al., 1996). However,

understanding such articles requires "specific attention outside of general literacy capacities" (Ad Hoc Committee on Health Literacy, 1999), rendering them effectively inaccessible to most (lay) people. Ideally, healthcare providers would stay current on all medical evidence and share relevant findings with patients, but this is impractical due to the volume and growth of the evidence base (Bastian et al., 2010; Marshall et al., 2021).

LLMs may provide a means for lay readers to access such findings by automatically summarizing and simplifying texts into plain language (August et al., 2023; Shaib et al., 2023). Done successfully, this could allow patients to access the most up-to-date literature relevant to their healthcare. In turn, this may promote health literacy broadly by disseminating trustworthy information (Thielmann et al., 2023; Cheng et al., 2022). But given the inherent risks to personal health, the *factual correctness* of such outputs is paramount in this domain.

While Shaib et al. (2023) showed that GPT-3 infrequently introduced outright errors when simplifying RCT abstracts, inaccuracies are occasionally introduced; this ought to be addressed before wide adoption of such technology. Unfortunately, there is no standard evaluation benchmark for factuality on this important *medical evidence text simplification* task. Consequently, it is unknown whether and to what degree existing automatic factuality evaluation metrics align with human judgments.

We posit that focusing on critical elements in the RCT structure is key for factual medical evidence communication. This work presents FACT-PICO, an expert-constructed factuality benchmark for the plain language summarization of technical abstracts describing RCTs. FACTPICO is a fine-grained benchmark focused on key characteristics

<sup>\*</sup>Equal Contribution

<sup>&</sup>lt;sup>1</sup>We focus on abstracts, because they are always publicly accessible, and typically include the key results that would be of interest to individuals.

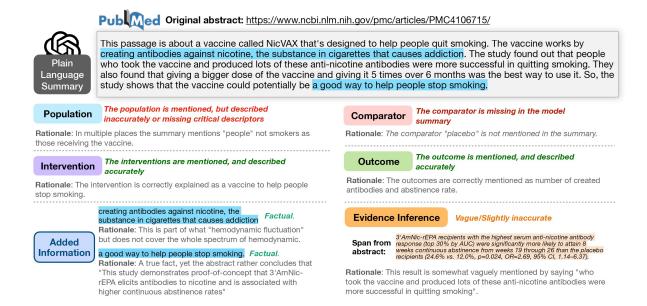


Figure 1: Expert evaluation of a GPT-4 plain language summary in FACTPICO. We omitted the original abstract (can be found in Appendix N) in this figure due to space limit. More examples in Appendix O.

of trials: Populations (e.g. *COVID patients; diabetics*), Interventions (e.g. *remdesivir*), Comparators (e.g. *placebo*), Outcomes (e.g. *30 day mortality, or pain*), as well as Evidence Inference (i.e., whether the intervention yielded a significant difference in the treated group with respect to the outcome; Lehman et al. 2019). PICO is a standard framework to structure clinical questions.

Figure 1 shows an example of FACTPICO annotation. In contrast to standard summaries, plain language summaries have the additional goal of simplifying content for lay readers. This may involve elaboration and explanation of difficult concepts to foster understanding (Srikanth and Li, 2021). Thus, FACTPICO includes a correctness assessment of added content. FACTPICO is also distinct in that it includes expert-written rationales that contextualize the evaluation of these fine-grained characteristics, providing a useful first step in assessing explainable factuality evaluation methods.

FACTPICO includes outputs from a mix of proprietary and open-source LLMs (GPT-4, Llama-2-Chat, and Alpaca). Our findings are somewhat less optimistic than prior work in medical (Shaib et al., 2023) and news summarization (Goyal et al., 2022): Factual errors (occasionally important ones) with respect to key RCT elements are introduced by LLMs. One concerning phenomenon is the extent to which models overgeneralize, resulting in problematic information loss (Trienes et al., 2024).

Using FACTPICO, we evaluated a suite of exist-

ing automatic metrics shown to perform well for factuality in summarization (Scialom et al., 2021; Goyal and Durrett, 2021; Fabbri et al., 2022; Zha et al., 2023a; Tang et al., 2023a), as well as newly devised LLM-based evaluations. We find that existing metrics correlate with expert ratings at the system level, but not at the instance level. The best performing metric is an LLM-based approach in which we first identify key RCT elements; this shows that providing models with explicit domain knowledge may help. Analysis of LLM-generated rationales shows that LLMs often provide flawed reasoning when justifying their self-evaluations.

#### 2 FACTPICO Benchmark

The FACTPICO benchmark consists of expert factuality assessments of 345 LLM-generated plain language summaries of 115 RCT abstracts.

## 2.1 LLM-based Plain Language Summaries

The medical abstracts used in FACTPICO are sourced from the *Evidence Inference 2.0* dataset (DeYoung et al., 2020), which contains abstracts and full articles that describe RCTs from PubMed. We use a subset of the abstracts which include annotated spans that state clinical results. We also exclude abstracts that have corresponding humanwritten plain language summaries in the PubMed database, which may have appeared in the pretraining data of the LLMs we are experimenting

with.<sup>2</sup> FACTPICO includes 115 abstracts randomly sampled from this filtered subset.

For each abstract, we generate plain language summaries using three different LLMs under zero-shot prompting. This zero-shot setting better emulates how lay users would likely interact with models, compared to few-shot. We use GPT-4 (OpenAI, 2023), Llama-2-Chat (Touvron et al., 2023), and Alpaca (Taori et al., 2023), which resulted in 345 (115  $\times$  3) total outputs. In a preliminary investigation, we found that these models seemed qualitatively best at generating plain language summaries. Details about this pilot study, prompt exploration, and the final prompts are in Appendix E.

#### 2.2 Human Evaluation Framework

We evaluate generated summaries using a set of questions addressing factuality as related to the key PICO aspects of RCTs, as well as information added by LLMs during simplification. We ask evaluators to **score** model outputs and provide **natural language rationales** to justify their ratings. Annotation interface details are in Appendix P.

PICO Elements. Population, Intervention, Comparator, Outcome (PICO) elements are the key components of an RCT (Richardson et al., 1995). The trial **population** concerns characteristics of subjects in the trial, including what condition they have, the number of participants, and their demographics. The **intervention** is the active treatment being assessed for comparative efficacy; the **comparator** is the control with respect to which this is being compared. Finally, **outcomes** are those things that are measured to determine results.

Accurate representation of these essential descriptors in the plain language summaries of an RCT is imperative. More explanations and examples are provided in Appendix Table 7. We ask evaluators to provide a rating between 1 to 4 codifying the factuality with respect to each PICO element, respectively: **4:** *Mentioned and described accurately;* **3:** *Mentioned but described somewhat inaccurately or vaguely;* **2:** *Mentioned but described with severe inaccuracies and/or is missing critical descriptors;* **1:** *Missing.* We describe the rating system in greater detail in Appendix A.

**Rationales.** Evaluators are also asked to provide natural language *rationales* justifying their chosen

rating. Such rationales may reveal technical justifications for annotations which would be difficult for a layperson to assess (see Figure 2 for examples). More usefully, rationales can express a degree of uncertainty in evaluations. For example, consider the following rationale for an intervention being evaluated as accurately mentioned:

The interventions are clearly described as one group receiving morphine and one triamcinolone, yet their combination with bupivacaine is missing, but doesn't seem very important.

This rationale indicates that the choice to evaluate the intervention as accurate was not a clear-cut decision. By including rationales in FACTPICO, this complex decision process is documented. We compare LLM rationales and expert rationales in Section 6, and envision future work to dive deeply into explainable factuality measures.

**Evidence Inference.** We additionally evaluate the evidence inference aspect of LLM-generated plain language summaries, i.e., whether the results and findings concerning PICO are reported factually, as PICO covers what outcomes were mea*sured*, it does not account for the corresponding results (e.g., if an intervention was found to outperform a comparator). A challenge here is that most trials will report multiple results, any of which may be conveyed (un)factually in the summary. Thus, we collect factuality assessments at the level of individual results, using findings annotated in the source abstracts from the Evidence Inference 2.0 dataset as reference. Evaluators are asked to determine how well each particular inference is reflected in the plain language summary (4: accurate, 3: vague/slightly inaccurate, 2: inaccurate, 1: not *mentioned*), with a free-text rationale.

Added Information. Unlike traditional summarization where content addition is seen as extrinsic hallucination (Maynez et al., 2020), plain language summarization often requires the model to explain and elaborate complex concepts (Srikanth and Li, 2021). It is important to verify the correctness of content additions. We ask annotators to highlight addition spans, determine whether they are factual, and justify each rating with a free-text rationale.

#### 2.3 Annotation

FACTPICO was annotated by two senior students in their fifth year of medical school who are highly proficient in English. They are experienced with

<sup>&</sup>lt;sup>2</sup>We check this by querying the Entrez database system and filtering out abstracts which have other abstracts associated with their PubMed ID.

Type	$\kappa$	Туре	$\kappa$
Population	0.56	Evidence Inference	0.47
Intervention	0.80	Added Info	0.65
Comparator	0.73	- correctness	0.86
Outcome	0.60		

Table 1: Inter-evaluator agreement. For PICO evaluations (left) we show agreement on the subset of 15 *undiscussed* texts (see Appendix C). For other *undiscussed* questions (right), we show agreement on all 75 doubly annotated texts.

data annotation for text simplification and summarization tasks. To ensure high annotation quality, we conducted a training phase which involved annotating a set of summaries as pilot (excluded from FACTPICO). Next, we collected two sets of annotations on 75 summaries (25 from each model). For the first 60 of these summaries, we asked annotators to discuss their responses with respect to the PICO elements. These discussions were intended to resolve any conceptual differences related to PICO annotation (Appendix C). The remaining 15 summaries were evaluated individually. Compensation was 12€ per hour, which is standard for fixed-hour contracts with the annotators' institutions.

**Inter-Evaluator Agreement.** Table 1 reports agreement on the 15 held-out (undiscussed) subset for PICO elements, as well as agreement on the full set of 75 texts for the Evidence Inference and Added Information questions (all annotated without discussion). We report Randolph's kappa (Randolph, 2005), a free-marginal version of Fleiss' kappa.<sup>3</sup> With respect to added information, the agreement reported in Table 1 is at the sentence level (for each sentence, whether it is considered as added information). The correctness of added information is reported on the 48 spans both annotators identified as added information. Overall, evaluators showed moderate to high agreement (Artstein and Poesio, 2008). This range of kappa values are expected given the degree of subjectivity inherent to rating the severity of factual errors.

## 3 Evaluation and Analysis of Plain Language Summaries

Results of the factuality evaluation of the plain language summaries generated by GPT-4, Llama-2, and Alpaca in FACTPICO are presented in Table 3.

We report average ratings for each PICO category and Evidence Inference, as well as the number and percentage of non-factual added information spans identified by either of the evaluators. Since readability is a key goal of *plain language* summarization, we also report the change in Flesch-Kincaid Grade Level (Kincaid et al., 1975) between the abstract and the summary ( $\Delta$ FK), as well as the ROUGE-L score (Lin, 2004), which quantifies n-gram overlap between the generated plain language summary and the original abstract (i.e., high overlap may indicate low readability).

Factuality vs. Simplicity. We observe a clear trade-off between the simplicity of the generated text and its factuality. Plain language summaries generated by Alpaca are rated as more factual across all dimensions, including most of the automatic metrics (Table 5) discussed in Section 5. However, Alpaca is also the most extractive (i.e., heavily relying on deletion), with a near 0.5 ROUGE-L score and an advanced reading level. By contrast, GPT-4 and Llama-2 both simplify by rephrasing, with similar ROUGE-L and reading levels. Yet the plain language summaries generated by GPT-4 and Llama-2 are less factual, with a significant increase in the number of hallucinations (added non-factual information). Comparatively, Llama-2 produced the least factual summaries.

Can LLMs reliably convey *critical* RCT elements? The fine-grained framework in FACT-PICO exposes issues with LLM-generated plain language summaries for medical texts. As shown in Table 3, LLMs can explain the Interventions and Outcomes more accurately, while failing to do so for other information, such as the Populations, Comparators, and Evidence Inference.

Overgeneralization and omissions of such critical elements may distort the findings and conclusions of medical research studies in nuanced but important ways. Consider this example where GPT-4 omitted a critical element of RCT in its summary.

Scientists did a study to see if people who know they have a gene that puts them at risk for heart problems would eat healthier. They looked at people's intake of a vitamin called folate, which is good for heart health. The participants in the study were told if they have the risk gene or not. But in the end, there was no difference in how much folate the people with the risk gene ate compared to those without the risk gene. So, knowing if they have the risk gene didn't make people eat healthier.

In the original RCT abstract, participants received "either general healthy eating advice or varying lev-

<sup>&</sup>lt;sup>3</sup>Free-marginal refers to distributions where raters have no prior knowledge as to the quantity of instances that would be assigned each rating, such as is the case with FACTPICO.

#### **Supplementary Elaboration (factual)**

Source: We conducted a randomized, placebo-controlled clinical trial to prove efficacy of botulinum toxin for neck pain in chronic whiplash syndrome. Simplification [GPT-4]: Scientists did a study to see if a certain type of medicine, called botulinum toxin, could help people who had long-lasting neck pain from whiplash (an injury that happens if your head jerks hard).

Rationale: This is an explanation that is not mentioned in the original.

#### **Substituting Elaboration (factual)**

Source: The aim of the study was to compare the effect of co-administered dexamethasone and diclofenac potassium (diclofenac K) with diclofenac K alone on the postoperative pain, swelling and trismus after surgical removal of third molars.

Simplification [Llama-2]: The study compared how well two different medicines, dexamethasone and diclofenac K, worked together to help people feel better after having surgery to remove their wisdom teeth

Rationale: This is the correct explanation for 3rd molar.

#### Nonsensical Hallucination (non-factual)

Source: Intra-articular morphine was as effective as intra-articular triamcinolone for analgesia in patients with osteoarthritis knee. Intra-articular morphine is possibly a better option than intra-articular steroid as it has lesser side effects.

Simplification [GPT-4]: The study concluded that both morphine and triamcinolone are useful, but morphine might be better as it has fewer side effects.

So, if granny or grandpa is having pain because of osteoarthritis, an injection with morphine could help them better.

Rationale: This is definitely not factual or good advice to give, even though the information within is technically true.

#### Jumping to Conclusions (non-factual)

Source: This study aimed to determine if raloxifene prevents the negative effects of diabetes on skeletal fragility in diabetes-prone rats.

Simplification [Llama-2]: Raloxifene is a medicine that can help prevent diabetes from causing weak bones

Rationale: The abstract says that this study aims to test IF Raloxifene can help prevent diabetes from causing "weak bones"/"if raloxifene prevents the negative effects of diabetes on skeletal fragility in diabetes-prone rats"

Table 2: Examples of added information (underlined) found in plain language summaries within FACTPICO.

	Pop.	Inter.	Comp.	Out.	E.Inf.	Average	#tokens	$\Delta$ <b>FK</b> $\uparrow$	Rg-L↓	#N↓	% N↓
ALPACA	3.30	3.70	3.42	3.77	3.46	3.53	170.95	-0.61	0.479	8	7.0
GPT-4	3.12	3.52	3.20	3.56	3.25	3.33	162.91	2.87	0.146	53	31.3
LLAMA-2	2.71	3.40	2.70	3.41	2.80	3.00	116.27	2.92	0.136	57	38.3

Table 3: Human evaluation on a 1-4 scale ↑ (Section 2.2) on the factuality of PICO elements and Evidence Inferences. The average length of the original technical abstracts is 343.5 tokens.  $\Delta$ FK and Rg-L: Flesch Kincaid Grade Level difference and ROUGE-L between abstract and summary. #N: number of non-factual additions. %N: percentage of summaries in FACTPICO with at least one non-factual addition.

els of personalised nutrition advice". One of the levels of such personalized advice involved testing for a gene "in relation to cardiovascular health and the importance of a sufficient intake of folate". But, the generated summary above fails to mention the comparator, which is the general healthy eating advice; instead, it focuses on the overgeneralized intervention, i.e., personalized advice involving genetic variants (without mentioning the personalized nutrition advice). These omissions not only render the simplification incomplete, but they also make understanding the actual result impossible.

# Do LLMs generate accurate elaborations? ble 3 also shows a concerning amount of nonfactual additions within the generated simplifications. Examples of these span-level annotations can be found in Table 2. Most of these non-factual additions may misrepresent the original medical text and consequently mislead lay readers. Even GPT-4 produced many plain language summaries (31.3% of all additions added) with such errors, raising questions as to the trustworthiness of LLMs for tasks in high-stakes domains such as medicine, when used by zero-shot prompting (a most common use-case for lay users).

## **Factuality Evaluation Metrics**

FACTPICO is a dataset that can be used to assess automatic evaluation methods for plain language summarization of RCT texts. We assess existing methods used for factuality evaluation, as well as the capabilities of LLMs themselves to evaluate factuality. These analyses focus on questions about PICO elements and Evidence Inference; we leave evaluation of the factuality of added information for future work, as this entails fact-checking using external knowledge sources.

#### 4.1 Factuality Metrics Evaluated

We first evaluate a suite of existing automatic factuality metrics shown to be effective in prior work. (1) Dependency-Arc Entailment (DAE) (Goyal and Durrett, 2021) decomposes summaries into smaller entailment tasks at the arc-level to assess their factuality. We use a numeric score by taking the minimum of the probability scores assigned to individual arcs. (2) QuestEval (Scialom et al., 2021) uses a QA-based framework to analyze the factual faithfulness of a summary to the original text. This method scores summaries a 0 if there is no common token and a 1 for an exact match. (3) QAFactEval (Fabbri et al., 2022) is a QA-

		QAFact	Quest	AlignS	DAE	GPT-4	Llama-2	Alpaca	Mistral	Extract
	Pop.	0.237	0.226	0.201	0.116	0.265	0.042	0.084	-0.011	0.270
	Inter.	0.232	0.202	0.117	0.126	0.223	0.059	0.048	0.093	0.234
_	Comp.	0.176	0.177	0.140	0.123	0.341	0.036	0.073	0.015	0.387
$ au_b$	Out.	0.228	0.214	0.216	0.130	0.276	0.078	0.028	0.082	0.266
	Evd.Inf	0.248	0.221	0.242	0.075	0.405	0.056	0.053	0.021	0.372
	Avg.	0.289	0.290	0.244	0.152	0.475	0.055	0.081	0.047	0.474
	Pop.	0.311	0.300	0.269	0.155	0.333	0.055	0.107	-0.012	0.349
	Inter.	0.298	0.261	0.150	0.168	0.276	0.074	0.059	0.111	0.293
	Comp.	0.232	0.236	0.178	0.163	0.434	0.046	0.092	0.019	0.494
ho	Out.	0.292	0.272	0.278	0.165	0.340	0.099	0.035	0.098	0.332
	Evd.Inf	0.337	0.300	0.329	0.099	0.524	0.074	0.072	0.027	0.490
	Avg.	0.406	0.412	0.348	0.219	0.619	0.080	0.115	0.065	0.633
	Pop.	0.459	0.454	0.444	0.408	0.476	0.374	0.390	0.334	0.472
	Inter.	0.332	0.321	0.291	0.294	0.357	0.282	0.286	0.348	0.343
$acc_{eq}$	Comp.	0.378	0.379	0.364	0.357	0.461	0.328	0.345	0.331	0.468
	Out.	0.337	0.332	0.333	0.302	0.383	0.293	0.283	0.350	0.363
	Evd.Inf	0.501	0.490	0.499	0.425	0.586	0.412	0.408	0.367	0.558
	Avg.	0.611	0.611	0.589	0.545	0.676	0.482	0.481	0.402	0.686

Table 4: Kendall's  $\tau_b$ , Spearman's  $\rho$ , and pairwise accuracy  $acc_{eq}$  of systems to human evaluations. Note, each FACTPICO attribute is compared against the overall score produced by each system, which for LLM evaluators is the average rating across each element. We present an attribute-wise comparison for LLM evaluators in Appendix F.

	QAFact	Quest	AlignS	DAE	GPT-4	Llama-2	Alpaca	Mistral	Extract
ALPACA	3.680	0.547	0.884	0.654	3.608	3.375	0.934	3.619	3.277
GPT-4	1.976	0.415	0.683	0.317	3.528	3.225	0.942	3.297	2.891
LLAMA-2	1.894	0.412	0.610	0.379	3.152	3.128	0.920	3.447	2.614

Table 5: Average of systematic metrics per LLM for plain language summary.

based metric that combines various components from other factuality metrics and assigns scores based on the *LERC* score (Chen et al., 2020). This score, usually used for evaluating reading comprehension answers, ranges from 1 to 5, where 1 is a completely wrong answer and 5 is a perfect answer. (4) **AlignScore** (Zha et al., 2023b) is an alignment-based method for analyzing factual consistency. The final score assigned is the average of the maximum alignment probabilities between sentences from the summary to context chunks from the original abstract.

#### 4.2 LLM Evaluators

Prior work has also shown that LLMs themselves can be good evaluators for factuality in summarization (Luo et al., 2023; Wang et al., 2023; Tang et al., 2023a; Tian et al., 2023). FACTPICO judgments are finer-grained. Therefore, we prompt LLMs with instructions emulating the questions asked of human evaluators in FACTPICO. In addition to ratings, the LLMs also generate rationales for their scores, which we analyze in Section 6.

The implementation details and the prompts used for this task are in Appendix H. Prior work found that LLM evaluations may be biased in that a system may 'prefer' its own outputs (Liu et al., 2023), but we did not observe this here (Table 5).

**Full-Text Evaluation.** As input, we provide LLMs the full-text of a complex medical abstract and corresponding LLM-generated plain language summary. We instruct the evaluator LLM to find PICO elements in the text and evaluate them according to the provided criteria. In addition to PICO elements, we also evaluate Evidence Inference outputs; here the reference results span (annotated in prior work) is compared against the full text of the plain language summary. We evaluate four LLMs in this way: GPT-4 (OpenAI, 2023), Llama-2-Chat (7B) (Touvron et al., 2023), Alpaca (7B) (Taori et al., 2023), and Mistral (7B-Instruct-v0.1; Jiang et al. 2023).

**PICO-R Extraction Pipeline.** We also evaluate LLM scorers explicitly informed of the PICO elements and results inferred from evidence (PICO-R). We adopt this two-stage pipeline using GPT-4. We first extract PICO-R from both the original abstract and the plain language summary.<sup>4</sup> For Evidence

<sup>&</sup>lt;sup>4</sup>Initial experiments showed that GPT-4 can more accurately extract PICO elements compared to other LLMs and PICO-tagger models (Nye et al., 2020), especially for plain language summaries.

Inference, this extraction is only necessary from the summary. Next, extracted elements are passed to GPT-4 along with an evaluation prompt.

## 5 Factuality Metric Evaluation Results

To conduct a meta evaluation of the system factuality metrics, we compute the Kendall's  $\tau_b$  coefficient, Spearman correlation coefficient, and Pairwise Accuracy coefficient (Deutsch et al., 2023) for automatic vs. human evaluations for each of the PICO and evidence inference aspect. We posit that good metrics must correlate well with the most salient elements in a high-stake domain.

**Instance-level Results.** We present the results for the above metrics in Table 4.<sup>5</sup> Across most measures, with the exception of Kendall's  $\tau_b$ , the pipeline combining GPT-4 evaluation with PICO extraction has the highest correlations with the FACTPICO ratings. Thus, decomposing the original evaluation task into separate localization then evaluation steps yield better performance, indicating that LLMs benefit from breaking apart a complex tasks into a series of simpler steps completed separately.

FACTPICO ratings correlated more with dedicated factuality models than LLM evaluators with the exception of GPT-4. In fact, the smaller, open-source LLMs had barely any correlation at all with the FACTPICO ratings. These results show that there is a sizeable gap in an LLMs ability to evaluate generated text compared to its generation capabilities.

**System-level Results.** LLMs in Table 5, with the exception of Alpaca, produce similar results compared to the human average PICO-R ratings in FACTPICO in Table 3. Alpaca rates plain language summaries low, with averages under the minimum in the scale (1), and so has failed to follow the evaluation instructions. Although Alpaca-generated summaries have higher ratings, these aren't meaningful as are extractive.

**Visual Analysis.** For more visual analysis of the system factuality metrics vs human ratings, Figure 2 shows the top three performing systems vs. averaged human scores. The highest performing system, the PICO-R extraction pipeline, has a the most balanced linear pattern. On the other hand,

	GPT-4	Llama-2	Alpaca	Mistral	Extract
P	0.189	0.073	0.153 0.21	0.156	0.201
I	0.173	0.073	<b>0.210</b> <sub>0.11</sub>	0.169	0.205
$\mathbf{C}$	0.181	0.092	$0.126_{\ 0.11}$	0.173	0.215
O	0.148	0.075	$0.111_{\ 0.19}$	0.158	0.176
R	0.061	0.019	$0.045_{\ 0.60}$	0.045	0.060

Table 6: BERTScore (rescore baseline) F1 average for PICO and Evidence Inference (R). Note, the results for Alpaca omit rationales. We show the percentage of nonempty rationales next to the BERTScore for Alpaca. We caution comparing the results for Alpaca here against those of other LLMs because instances of empty rationales have been excluded in this evaluation.

GPT-4 without PICO-R extraction often rates non-factual summaries highly. Contrastively, QuestE-val, the best performing non-LLM factuality metric, is more cautious and rates high quality summaries lower. We further analyze automatic vs. human rating distributions in Appendix G. These results hint at the potential challenges of factuality assessment brought by the shift in readability.

## 6 Preliminary LLM Rationale Analysis

We perform a preliminary analysis on LLM rationales, comparing to expert rationales in FACT-PICO; we leave a thorough human evaluation of LLM rationales for future work.

Qualitative analysis of a small sample of rationales show that Llama-2 and Mistral are often able to comprehend the medical text but fail in making correct judgments according to the provided instructions. Most commonly, Mistral focuses on the abstract rather than the evaluated summary, while Llama-2 generates long explanations that eventually arrive at the wrong conclusion. Rationales from GPT-4 and its pipelined counterpart usually did not do this. For the most part, rationales made logical sense. However, some rationales were overly generous in it evaluation, ignoring critical errors. We provide examples of such erroneous rationales in Appendix M.

Table 6 shows the BERTScore (Zhang et al., 2020) between expert vs. LLM rationales across all evaluators experimented in this work. Overall, the rationales are dissimilar. The PICO-R extraction pipeline produced rationales most similar to that of humans, with the exception of rationales for evidence inferences, where GPT-4 rationales are slightly more similar than others. Despite explicit prompting, Alpaca frequently did not produce rationales and stopped generation after outputting a

<sup>&</sup>lt;sup>5</sup>The evaluation scale for Llama-2 was reversed because it struggled to follow the original instructions. We described this behavior in greater detail in Appendix K

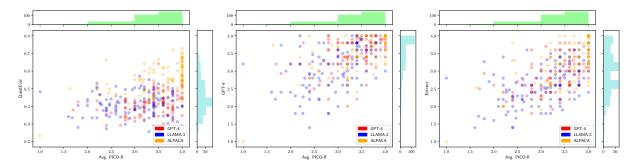


Figure 2: QuestEval (left), GPT-4 eval (mid), and Extract (right) against Avg. PICO-R (x-axis) for plain language summaries generated by GPT-4 (red), Llama-2 (blue), and Alpaca (orange). Label distributions shown on the sides.

numerical rating.

Despite reporting these preliminary results, we advise caution when adopting reference-based metrics for automatic free-text rationale evaluation, as the results may be misleading. For example, we observed that Alpaca rationales tend to be a somewhat arbitrary variation of the rating descriptions (e.g., The interventions were described accurately), and human evaluators often use the words in these descriptions in their rationales (e.g., The intervention is described correctly as an 8 week program of specific exercises), inflating word overlap scores. By contrast, an example Llama-2 rationale (The comparator in the PICO model is the placebo group...) could present the same idea as a human rationale (ex: Comparator (placebo) is mentioned.), but result in a low BERTScore (ex: -0.182) because rationales from Llama-2 tend to (much) longer than those written by humans. Future work developing metrics for rationale correctness should not solely rely on reference-based metrics. We additionally present a length analysis in Appendix L.

## 7 Related Work

Meta-evaluation of factual consistency metrics in summarization (and related tasks) in the "general domain", e.g., Wikipedia and news, has garnered considerable attention (Pagnoni et al., 2021; Honovich et al., 2022; Laban et al., 2022; Tang et al., 2023a; Min et al., 2023). However, in addition to focusing on different domains, these existing benchmarks in summarization include primarily older models. Newer LLMs may yield novel error types (or be more factual overall) (Tang et al., 2023a).

Existing summarization factuality benchmarks also fail to generalize to *simplification*, in which content addition in the form of elaborations or explanations is often necessary (Srikanth and Li,

2021). Broadly, simplification entails substantial language changes that often lead to the text being more general (Li and Nenkova, 2015). Devaraj et al. (2022) evaluated the factuality of automated simplification model outputs at the sentence-level, noticing that content *deletion* can often lead to factual errors (in contrast to only summarization). Our findings confirm that these errors also exist in *plain language* summarization; the overgeneralization problem may lead to safety issues in the medical domain.

Shaib et al. (2023) and Tang et al. (2023b) evaluated LLM-generated summaries of medical evidence. Notably, Shaib et al. (2023)'s work included an evaluation of *plain language* summaries. Our work deepens this analysis with a finer-grained evaluation focusing on critical components of RCTs and medical evidence, covering three LLMs. Our findings call for caution against LLM-generated plain language summaries despite the absence of outright inconsistencies. Pal et al. (2023) presents an analysis of "hallucinations" in medical QA tests, focusing on reasoning rather than factual consistency.

The inclusion of human-written natural language rationales in factuality benchmarks is rare, and there is a paucity of work evaluating these. The FELM benchmark (Chen et al., 2023), an opendomain evaluation of LLM-generated long form texts covering factual knowledge, math, and reasoning included human rationales. Work in LLM critiquing has started to incorporate natural language critiques from both humans (Saunders et al., 2022) and LLMs (Cui et al., 2023; Kim et al., 2024). FACTPICO is the first factuality benchmark of the plain language summarization task that includes expert-generated natural language rationales.

#### 8 Conclusions

We introduced FACTPICO, an expert-annotated benchmark in the domain of evidence-based medicine for evaluating the factuality of plain language summarization with respect to clinically important dimensions. Using FACTPICO, we presented an analysis of factual errors along these fine-grained aspects in LLM-generated plain language summaries. We also presented an analysis of methods to evaluate factuality, including both dedicated factuality models and novel LLM-based methods.

## Limitations

The process of evaluating texts in FACTPICO was time consuming for human evaluators, requiring close reading of complex, technical language. Consequently, to make the workload manageable, we did not ask evaluators to localize PICO elements (or results regarding these) in texts in plain language summaries. Annotating how these elements are represented overall in the summary would provide more insights in cases where they are vaguely represented and thus required more focus in this evaluation. We automatically localize these elements in the evaluated PICO-R extraction pipeline, but we encourage future work expanding FACTPICO to include human-annotated, span-level annotations.

We primarily evaluated the zero-shot capabilities of LLMs in evaluating the factuality of simplified medical texts. We chose this setting as it best reflects how an end-user—a lay individual—would likely interact with an LLM, as such users are unlikely to provide LLMs with expert-evaluated plain language summaries for few-shot prompting. Future work could explore other LLM evaluation methods methods and use FACTPICO as an evaluation benchmark.

In our meta-evaluation we compared factuality metrics that assess the overall factuality of a text against numerical ratings assessing the fine-grained factuality of key characteristics in these texts. We acknowledge that this is not an equivalent comparison. However, we posit that the aggregation of these assessments should be well-correlated with the overall factuality of RCT texts for them to be useful in this important domain.

## Acknowledgments

We thank David Heineman for help with the annotation interface, Yao Dou for help with Alpaca,

and Tanya Goyal and Greg Durrett for useful discussions. This research is partially supported by NSF CAREER Award IIS-2145479 and Good Systems, 6 a UT Austin Grand Challenge to develop responsible AI technologies. Trienes, Göke and Coers were supported by the Cancer Research Center Cologne Essen (CCCE). Trienes was also supported by the Federal Ministry of Education and Research (BMBF) and by a fellowship within the IFI programme of the German Academic Exchange Service (DAAD). Wallace was supported in this work by the National Institutes of Health (NIH), grant R01LM012086, and by the National Science Foundation (NSF), grant 1750978. Xu is supported in part by NSF awards IIS-2144493 and IIS-2112633.

#### References

Ad Hoc Committee on Health Literacy. 1999. Health literacy: report of the council on scientific affairs. *Journal of the American Medical Association*, 281(6):552–557.

Ron Artstein and Massimo Poesio. 2008. Inter-coder agreement for computational linguistics. *Computational linguistics*, 34(4):555–596.

Tal August, Lucy Lu Wang, Jonathan Bragg, Marti A Hearst, Andrew Head, and Kyle Lo. 2023. Paper plain: Making medical research papers approachable to healthcare consumers with natural language processing. *ACM Transactions on Computer-Human Interaction*, 30(5):1–38.

Hilda Bastian, Paul Glasziou, and Iain Chalmers. 2010. Seventy-five trials and eleven systematic reviews a day: how will we ever keep up? *PLoS medicine*, 7(9):e1000326.

Anthony Chen, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2020. MOCHA: A dataset for training and evaluating generative reading comprehension metrics. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6521–6532, Online. Association for Computational Linguistics.

Shiqi Chen, Yiran Zhao, Jinghan Zhang, I-Chun Chern, Siyang Gao, Pengfei Liu, and Junxian He. 2023. Felm: Benchmarking factuality evaluation of large language models. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Christina Cheng, Emma Gearon, Melanie Hawkins, Crystal McPhee, Lisa Hanna, Roy Batterham, and Richard H Osborne. 2022. Digital health literacy as a

<sup>&</sup>lt;sup>6</sup>https://goodsystems.utexas.edu

- predictor of awareness, engagement, and use of a national web-based personal health record: Population-based survey study. *Journal of Medical Internet Research*, 24(9):e35772.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. 2023. Ultrafeedback: Boosting language models with high-quality feedback. *arXiv* preprint arXiv:2310.01377.
- Daniel Deutsch, George Foster, and Markus Freitag. 2023. Ties matter: Meta-evaluating modern metrics with pairwise accuracy and tie calibration. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12914–12929, Singapore. Association for Computational Linguistics.
- Ashwin Devaraj, William Sheffield, Byron Wallace, and Junyi Jessy Li. 2022. Evaluating factuality in text simplification. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7331–7345, Dublin, Ireland. Association for Computational Linguistics.
- Jay De Young, Eric Lehman, Benjamin Nye, Iain Marshall, and Byron C. Wallace. 2020. Evidence inference 2.0: More data, better models. In *Proceedings of the 19th SIGBioMed Workshop on Biomedical Language Processing*, pages 123–132, Online. Association for Computational Linguistics.
- Alexander Fabbri, Chien-Sheng Wu, Wenhao Liu, and Caiming Xiong. 2022. QAFactEval: Improved QAbased factual consistency evaluation for summarization. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2587–2601, Seattle, United States. Association for Computational Linguistics.
- Tanya Goyal and Greg Durrett. 2021. Annotating and modeling fine-grained factuality in summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1449–1462, Online. Association for Computational Linguistics.
- Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022. News summarization and evaluation in the era of GPT-3. *arXiv preprint arXiv:2209.12356*.
- Eduardo Hariton and Joseph J Locascio. 2018. Randomised controlled trials—the gold standard for effectiveness research. *BJOG: an international journal of obstetrics and gynaecology*, 125(13):1716.
- David Heineman, Yao Dou, and Wei Xu. 2023. Thresh: A unified, customizable and deployable platform for fine-grained text evaluation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 336–345, Singapore. Association for Computational Linguistics.

- Or Honovich, Roee Aharoni, Jonathan Herzig, Hagai Taitelbaum, Doron Kukliansy, Vered Cohen, Thomas Scialom, Idan Szpektor, Avinatan Hassidim, and Yossi Matias. 2022. TRUE: Re-evaluating factual consistency evaluation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3905–3920, Seattle, United States. Association for Computational Linguistics.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b.
- Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoo Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2024. Prometheus: Inducing finegrained evaluation capability in language models. In *The Twelfth International Conference on Learning Representations*.
- J Peter Kincaid, Robert P. Fishburne Jr., Richard L. Rogers, and Brad S. Chissom. 1975. Derivation of new readability formulas (Automated Readability Index, Fog Count and Flesch Reading Ease Formula) for navy enlisted personnel. *Naval Technical Training Command Millington TN Research Branch*.
- Philippe Laban, Tobias Schnabel, Paul N. Bennett, and Marti A. Hearst. 2022. SummaC: Re-visiting NLI-based models for inconsistency detection in summarization. *Transactions of the Association for Computational Linguistics*, 10:163–177.
- Eric Lehman, Jay DeYoung, Regina Barzilay, and Byron C. Wallace. 2019. Inferring which medical treatments work from reports of clinical trials. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3705–3717, Minneapolis, Minnesota. Association for Computational Linguistics.
- Junyi Jessy Li and Ani Nenkova. 2015. Fast and accurate prediction of sentence specificity. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, AAAI'15, page 2281–2287. AAAI Press.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-eval: NLG evaluation using gpt-4 with better human alignment. In *Proceedings of the 2023 Conference on*

- Empirical Methods in Natural Language Processing, pages 2511–2522, Singapore. Association for Computational Linguistics.
- Zheheng Luo, Qianqian Xie, and Sophia Ananiadou. 2023. Chatgpt as a factual inconsistency evaluator for text summarization.
- Iain James Marshall, Veline L'Esperance, Rachel Marshall, James Thomas, Anna Noel-Storr, Frank Soboczenski, Benjamin Nye, Ani Nenkova, and Byron C Wallace. 2021. State of the evidence: a survey of global disparities in clinical trials. *BMJ Global Health*, 6(1):e004145.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, Online. Association for Computational Linguistics.
- Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–12100, Singapore. Association for Computational Linguistics.
- Benjamin Nye, Ani Nenkova, Iain Marshall, and Byron C. Wallace. 2020. Trialstreamer: Mapping and browsing medical evidence in real-time. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 63–69, Online. Association for Computational Linguistics.
- OpenAI. 2023. Gpt-4 technical report.
- Artidoro Pagnoni, Vidhisha Balachandran, and Yulia Tsvetkov. 2021. Understanding factuality in abstractive summarization with FRANK: A benchmark for factuality metrics. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4812–4829, Online. Association for Computational Linguistics.
- Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2023. Med-halt: Medical domain hallucination test for large language models. In *Proceedings of the 27th Conference on Computational Natural Language Learning (CoNLL)*, pages 314–334.
- Justus J Randolph. 2005. Free-marginal multirater kappa (multirater k [free]): An alternative to fleiss' fixed-marginal multirater kappa. *Online submission*.
- W S Richardson, M C Wilson, J Nishikawa, and R S Hayward. 1995. The well-built clinical question: a key to evidence-based decisions. *ACP J. Club*, 123(3):A12–3.

- David L Sackett, William MC Rosenberg, JA Muir Gray, R Brian Haynes, and W Scott Richardson. 1996. Evidence based medicine. *BMJ: British Medical Jour*nal, 313(7050):170.
- William Saunders, Catherine Yeh, Jeff Wu, Steven Bills, Long Ouyang, Jonathan Ward, and Jan Leike. 2022. Self-critiquing models for assisting human evaluators. *arXiv preprint arXiv:2206.05802*.
- Thomas Scialom, Paul-Alexis Dray, Patrick Gallinari, Sylvain Lamprier, Benjamin Piwowarski, Jacopo Staiano, and Alex Wang. 2021. Questeval: Summarization asks for fact-based evaluation.
- Chantal Shaib, Millicent Li, Sebastian Joseph, Iain Marshall, Junyi Jessy Li, and Byron Wallace. 2023. Summarizing, simplifying, and synthesizing medical evidence using GPT-3 (with varying success). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1387–1407, Toronto, Canada. Association for Computational Linguistics.
- Neha Srikanth and Junyi Jessy Li. 2021. Elaborative simplification: Content addition and explanation generation in text simplification. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 5123–5137, Online. Association for Computational Linguistics.
- Liyan Tang, Tanya Goyal, Alex Fabbri, Philippe Laban, Jiacheng Xu, Semih Yavuz, Wojciech Kryscinski, Justin Rousseau, and Greg Durrett. 2023a. Understanding factual errors in summarization: Errors, summarizers, datasets, error detectors. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11626–11644, Toronto, Canada. Association for Computational Linguistics.
- Liyan Tang, Zhaoyi Sun, Betina Idnay, Jordan G Nestor, Ali Soroush, Pierre A Elias, Ziyang Xu, Ying Ding, Greg Durrett, Justin F Rousseau, Chunhua Weng, and Yifan Peng. 2023b. Evaluating large language models on medical evidence summarization. *npj Digital Medicine*, 6(1):158.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford\_alpaca.
- Rosa R L C Thielmann, Ciska Hoving, Jochen W L Cals, and Rik Crutzen. 2023. The effects of online access to general practice medical records perceived by patients: Longitudinal survey study. *Journal of Medical Internet Research*, 25:e47659.
- Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher Manning, and Chelsea Finn. 2023. Fine-tuning language models for factuality. In *NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and finetuned chat models.

Jan Trienes, Sebastian Joseph, Jörg Schlötterer, Christin Seifert, Kyle Lo, Wei Xu, Byron C. Wallace, and Junyi Jessy Li. 2024. InfoLossQA: Characterizing and recovering information loss in text simplification. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.

Jiaan Wang, Yunlong Liang, Fandong Meng, Zengkui Sun, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. 2023. Is ChatGPT a good NLG evaluator? A preliminary study. In *Proceedings of the 4th New Frontiers in Summarization Workshop*, pages 1–11, Singapore. Association for Computational Linguistics.

Yuheng Zha, Yichi Yang, Ruichen Li, and Zhiting Hu. 2023a. AlignScore: Evaluating factual consistency with a unified alignment function. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11328–11348, Toronto, Canada. Association for Computational Linguistics.

Yuheng Zha, Yichi Yang, Ruichen Li, and Zhiting Hu. 2023b. AlignScore: Evaluating factual consistency with a unified alignment function. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11328–11348, Toronto, Canada. Association for Computational Linguistics.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

#### **A Evaluation Guidelines**

Annotators are shown a "source" input (abstract of a technical paper describing an RCT), along with a plain language summary of it, automatically produced by a model, and asked a series of 7 questions for evaluation.

#### **Question I: Added Information**

For this part of the question, please highlight words, phrases, or sentences in the output that adds or modifies the information from the original document. Afterwards, answer whether this added/modified information is factual or not factual. Also, provide a rationale for why this is added/modified information and why it is factual or not factual.

# Questions II-V (one analogous question for each PICO element):

For instance, for **population**, we ask:

Ш	The population is mentioned, and described accurately
	The population is mentioned, but described somewhat inaccurately or vaguely
	The population is mentioned, but described with severe inaccuracies and/or is missing critical descriptors
	The population is missing in the model summary
	N/A

Please provide a rationale for why you chose the answer choice in relation to the abstract and the summary.

Some examples will actually not quite be randomized trials (e.g., they might be observational studies, or a description of a prospective trial not yet run). In these cases, it may be that there is no meaningful population (or intervention, comparator, and outcome). Here you should select the "N/A" option.

Exhaustive Outcomes (Analyzed in Appendix I): Sometimes, the plain language summary may not mention all the outcome measures described in the abstract. The summary may still be considered factual if the omitted measures are non-critical for the experiment and are not mentioned any further in the abstract. However, please do annotate

separately when the plain language summary does exhaustively mention all outcome measures and when it does not.

## **Question VI: Evidence Inferences**

You will be presented with a span of text highlighting the inferred result from the experiment presented in the abstract. Based on this span, choose the following based on how this span is presented in the plain language summary:

Accurate
Vague/Slightly Inaccurate
Inaccurate
Not mentioned

Please provide a rationale for why you chose the answer choice in relation to the evidence inference span and the summary.

## A.1 Additional Questions

We also collected information for any additional comments on the generated plain language summary, as well as contradictions that are *not* covered by the other questions. These annotations are scarce, thus we have not included them in the FACTPICO benchmark.

Additional comments: "asks you to write down any commands you would want the machine to follow if you could interact with it, e.g., "Make it shorter", "Explain XXX a bit more", and so on."

Contradictions: "Here you will be looking for content in the output that contradicts some part of the input. We ask you to annotate both the input and the output for this question. Please provide a rationale for why the content is a contradiction." The contradictions are analyzed in Appendix J.

## **B** Data Release and License

We reused RCT abstracts from the *Evidence Inference V2.0* dataset (DeYoung et al., 2020); https://evidence-inference.ebm-nlp.com, accessed 2024-02-15). All articles in this dataset are from the PubMed Open Access subset which only includes license terms that allow reuse (https://ncbi.nlm.nih.gov/pmc/tools/openftlist, accessed 2024-02-15). After discussion with our institutions' librarian on fair use, we release the annotations in FACTPICO under CC-BY-4.0.

PICO Element	Description	Critical Descriptors	Example(s)
Population	The types of patients involved in the trial	Demographics     Specific condition	<ul><li>Diabetic males</li><li>Healthy adults</li></ul>
Intervention	The treatments considered	<ul><li> Mentioned as an intervention</li><li> Differentiates from comparator</li></ul>	Aspirin
Comparator	The alternative treatment to which the intervention is being compared against.	<ul> <li>Mentioned as comparing against intervention</li> <li>Differentiates from intervention</li> </ul>	• Placebo
Outcome	The measures used. NOT what was found in the study ("result"). For example, if the study finds a drug reduces the duration of headache, the outcome here is just the "duration of headache", not that it reduced it.	Primary outcomes mentioned (cannot make any conclusions without them).	Mortality     Duration of headache

Table 7: PICO elements. A critical descriptor is a characteristic that is absolutely crucial to understanding the study.

Туре	$\kappa$
Population	0.47/0.56
Intervention	0.59/0.80
Comparator	0.63/0.73
Outcome	0.56/0.60

Table 8: Inter-evaluator agreement measured through Randolph's  $\kappa$ ; for PICO evaluations we show agreement on all 75 doubly annotated documents (left) and only the subset of 15 undiscussed documents (right).

# C Improving Agreement Through Discussion

Initially, we had observed low agreement on PICO questions among evaluators on the first 30 doubly annotated documents. Upon deeper analysis of these disagreements, we modified our instructions to be clearer and more detailed. We had asked evaluators to independently re-evaluate the previous annotations as well as doubly annotate 30 more documents. We observed that while agreement had improved significantly overall, evaluators still disagreed substantially on questions regarding population and outcome.

To fix this issue, we facilitated a "soft" discussion between evaluators regarding their annotations. Evaluators were presented with documents in which they had disagreed majorly on any PICO questions from the last 30 documents they evaluated. Then they were asked to come to consensus on how to rate these questions. Afterwards, we asked evaluators to independently reevaluate their previous annotations based on the insights they gained from this discussion. Similarly, they were asked to doubly annotate 15 new documents inde-

pendently.

Table 8 shows the resulting agreement through Randolph's kappa after this discussion. For each question type, the first number is the kappa for the entire set of 75 doubly annotated documents. The second number is the kappa for the set of 15 documents that were annotated independently after the discussion. Both sets of number show moderate to high agreement for these questions. Furthermore, the agreement for this undiscussed set being significantly higher than the agreement for all 75 doubly annotated documents indicates that this discussion method was effective at improving agreement.

## **D** Model Details and Compute

We used a High-RAM T4 GPU through Google Colab Pro+ to conduct our experiments.

## **D.1** Plain Language Summary Generation

**GPT-4** We used a frequency penalty of 0, presence penalty of 0, temperature of 1, and top p of 1.

**Llama-2 7B Chat**. We set the max new tokens to 4000, did multinomial sampling, temperature of 1, top k of 50, and top p of 1.0.

**Alpaca (7B)**. We set the max new tokens to 4000, used greedy decoding, temperature of 1, top k of 50, and top p of 1.0.

## **D.2** LLM evaluation

**GPT-4** We used a frequency penalty of 0, presence penalty of 0, temperature of 1, and top p of 1.

**Together.AI** We ran experiments with Llama-2 7B Chat, Alpaca (7B), and

Mistral-7B-Instruct-v0.1 using the Together.AI API interface. For all models, we set the max new tokens to 256, temperature of 0.6, top k of 90, top p of 0.8, and a repetition penalty of 1.1.

## E Zero-shot Prompts for Plain Language Summarization/Simplification

# E.1 Preliminary Prompt Exploration and Model Evaluation

We conducted a preliminary evaluation in Table 9 for prompt engineering and model exploration on 100 medical abstracts from Shaib et al. (2023). The prompts we tested for are all *typical-use* prompts, representing how these systems would be usually used by the public. We avoided engineering prompts for the best-case use of these systems for this task, as we did not want to induce a false trust in any of these systems for medical use, which could be potentially harmful. Given this criteria, we evaluated the following prompts:

## Paper Plain: (August et al., 2023)

My fifth grader asked me what this passage means: [abstract] I rephrased it for him, in plain language a fifth grader can understand:

## **Short**:

"My fifth grader asked me what this passage means: [abstract] Help me summarize it for him, in plain language a fifth grader can understand. Make it short."

#### **Summarize:**

"My fifth grader asked me what this passage means: [abstract] Help me summarize it for him, in plain language a fifth grader can understand."

#### 5th grade:

5th Grade: Paraphrase this passage completely in your own words. Always define words the reader may not know: [abstract]

### **Complex:**

"Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

### ### Instruction:

Rewrite the following complex passage in order to make it easier to understand by non-native

Model <sub>Prompt</sub>	DA	FK	#tokens
GPT-4 <sub>Paper Plain</sub>	85.93	9.155	216.77
GPT-4 <sub>5th grade</sub>	89.90	10.606	308.73
GPT-4 <sub>Summarize</sub>	84.90	9.583	183.71
GPT-4 <sub>Short</sub>	85.21	14.741	111.28
Flan-T5 <sub>Plain</sub>	87.53	14.741	47.43
Flan-T5 <sub>5th grade</sub>	81.44	15.031	28.24
LLAMA-2 <sub>Paper Plain</sub>	81.03	8.218	135.38
ALPACA <sub>Complex</sub>	88.41	13.308	113.21
ALPACA <sub>Medical</sub>	89.28	13.523	101.37
Dataset	-	11.879	293.74

Table 9: ChatGPT-DA (Wang et al., 2023), Flesch-Kincaid Grade Level, and # of tokens with spaCy tokenizer for preliminary evaluation. The Model<sub>Prompt</sub> we chose for FACTPICO simplifications are in green.

speakers of English.

### Input:"
[abstract]

###Response:"

#### **Medical:**

"Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

#### ### Instruction:

Rewrite the following medical abstract in order to make it easier to understand by non-native speakers of English.

### Input:"
[abstract]

###Response:"

## **E.2** FACTPICO **Prompts**

For GPT-4, we randomly sampled from GPT- $4_{Summarize}$  and GPT- $4_{Short}$ . We also used Llama- $2_{Paper\ Plain}$  and Alpaca<sub>Complex</sub>.

### **E.3** Explanation of Selection Criteria

We primarily selected prompts based on the length of the plain language summaries produced when they are used. They had to be substantially shorter than their corresponding abstracts, qualifying them as plain language *summaries*. We did not use any form of length control outside of prompt instructions, as we found prompt instructions, such as "Make it short.", are just as effective as other forms

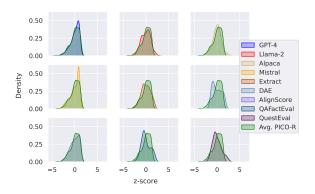


Figure 3: Plots of estimated Gaussian probability density functions from the standardized distributions of evaluated metrics.

of length control without creating any unwanted artifacts and disfluencies in the text. We did not consider ChatGPT-DA (Wang et al., 2023) to be a reliable measure of factuality. However, it was an effective sanity check in determining the relevance of the summaries to their corresponding abstracts.

We included two prompts to get generations from GPT-4 because both prompts met the above selection criteria and added more diversity in FACT-PICO. As a benchmark for factuality evaluation methods, including more diverse outputs will ensure only very robust systems are capable of attaining good results on FACTPICO.

## F Analysis of PICO Specific LLM Evaluations

We present the attribute-wise results for each LLMbased system we evaluated in Table 10. Similar to the results presented in Table 4, the GPT-4 based systems are more correlated with human ratings than the other LLM evaluators. Interestingly, we do see that for some attributes, in particular Population and Evidence Inference, the full text GPT-4 approach produced more human correlated ratings compared to those produced using the extraction pipeline. This suggests that evaluation of these elements requires more context than that provided in the extracted spans to perform a more accurate evaluation. However, despite this result, in terms of overall factuality, the extraction pipeline is slightly more correlated than the full text approach, suggesting an overall benefit to decomposing a complex task such as this into smaller, simpler steps to be executed sequentially.

## **G** System to Avg. PICO-R Visualizations

Figure 4 shows all system evaluations plotted against human evaluations. Figure 3 describes the distribution of averaged FACTPICO ratings alongside those of the automatic evaluation methods. These graphs display Gaussian approximations of the standardized distributions, allowing for easy visual comparison despite differences in scales. We only observe the average of the PICO-R ratings across human and LLM evaluations.

Our primary focus in this analysis is observing how closely the distribution of automaticallyderived factuality scores are to the human evaluations in FACTPICO. Upon visual inspection of Figure 3, we see that GPT-4, Mistral, and Align-Score seem to have the most closely aligned distributions to the average human rating distribution in FACTPICO. In order to quantify this, we also calculated the KL divergence between the distributions of the evaluated metrics and the distribution of the averaged human ratings from FACTPICO as they are presented in Figure 3. These results can be found in Table 11. As previously hypothesized, the three distributions closest to the averaged human distribution are from AlignScore, Mistral, and GPT-4, in that order. Another trend that is observable in Figure 3 are the differences in the narrowness of the distributions between some LLM evaluators and dedicated factuality models. Ratings from GPT-4, Mistral, and, to a certain extent, Alpaca typically have narrow distributions, indicating an overall preference towards a single rating. On the other hand, distributions from dedicated factuality models are wider, signaling more variation.

Interestingly, these is a mismatch between wellcorrelated metrics and metrics whose distribution aligns closer to that of the human ratings (Figures 2 and 4). Having a closely aligned distribution does not imply good correlation. Similarly, good correlation, unless it is exceptionally high, does not imply closely aligned distributions. This is the most evident in the case of AlignScore. However, this mismatch does highlight interesting trends. A left or right skew in a score distribution references the "strictness" of the evaluator, with a more leaning skew indicating the evaluators are less strict and vice versa. Through this lens, the human evaluators could be viewed as not as "strict" compared to many of the metrics. A possible explanation of this effect could be that the simplified nature of the text or the accurate elaborations present in the text

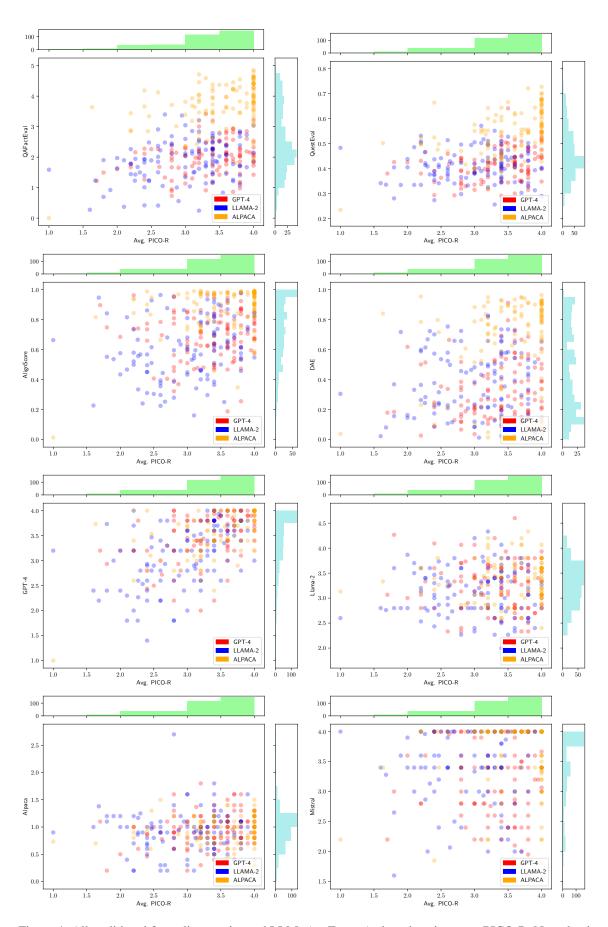


Figure 4: All traditional factuality metrics and LLMs (no Extract) plotted against avg. PICO-R. Note that human and LLMs scores are flipped (as 5-original) to be consistent with metrics in Section 4.1, hence higher is better.

		GPT-4	Llama-2	Alpaca	Mistral	Extract
	Pop.	0.323	0.017	0.095	-0.052	0.377
	Inter.	0.233	0.053	0.036	0.081	0.222
_	Comp.	0.506	0.001	0.064	-0.001	0.575
$ au_b$	Out.	0.260	0.086	0.090	0.145	0.261
	Evd. Inf.	0.561	0.137	-0.045	0.068	0.506
	Avg.	0.475	0.055	0.081	0.047	0.474
	Pop.	0.363	0.019	0.110	-0.059	0.429
	Inter.	0.246	0.058	0.040	0.089	0.246
	Comp.	0.558	0.001	0.072	-0.001	0.649
ho	Out.	0.279	0.094	0.100	0.159	0.292
	Evd. Inf.	0.667	0.170	-0.051	0.078	0.604
	Avg.	0.619	0.080	0.115	0.065	0.633
	Pop.	0.469	0.331	0.365	0.287	0.482
	Inter.	0.538	0.379	0.365	0.420	0.419
0.00	Comp.	0.604	0.339	0.366	0.343	0.632
$acc_{eq}$	Out.	0.502	0.381	0.383	0.450	0.423
	Evd. Inf.	0.665	0.419	0.342	0.287	0.631
	Avg.	0.676	0.482	0.481	0.402	0.686

Table 10: Kendall's  $\tau_b$ , Spearman's  $\rho$ , and pairwise accuracy  $acc_{eq}$  of PICO specific systems to human evaluations.

Method	KL
GPT-4	0.016
GPT-4 + PICO Extract	0.111
LLAMA-2	0.161
ALPACA	0.192
Mistral	0.042
DAE	0.238
AlignScore	0.015
QAFactEval	0.283
QuestEval	0.406

Table 11: KL divergence between the standardized distributions of evaluated metrics and that of the averaged PICO rating in FACTPICO

could have been confused for factual errors.

## H Zero-shot Prompts for LLM Evaluation

Here we present the prompts used for LLM evaluation. To find the implementation details see Appendix D.2.

## **H.1 Post-processing Ratings**

The ratings produced from the prompts displayed below follow a reverse scale, where lower scores indicate the evaluated text is more factual. For the sake of comparison, we post-process these ratings as 5 — Original Rating, flipping the scale so that it follows the rest of the evaluated metrics. Llama-2 is the only system that was not post-processed. We present the reasons for why this was done in Appendix K.

# **H.2** LLM Full-text Evaluation Prompt for PICO Elements

The following prompt is provided as a system prompt to the large language model.

You are given an abstract and a summary. <PI-COInfo> Find the <PICOElem> in accordance with PICO in both the abstract and the summary and use it rate the summary between 1 to 5.

The ratings are as follows.

- 1 The <PICOElem> is mentioned in the model summary and described accurately.
- 2 The <PICOElem> is mentioned in the model summary but described vaguely or somewhat inaccurately.
- 3 The <PICOElem> is mentioned in the model summary but described inaccurately or is missing critical descriptors.
- 4 The <PICOElem> is missing in the model summary.

5 - N/A

Please provide only the rating and the rationale for the rating. Provide the rating after stating "Rating:".

The model is then queried as follows:

Abstract: <Abstract>

passage. Pull direct quotes from the passage:

Summary: <Summary>

The tag "<PICOInfo>" correspond to the following four descriptors depending on the evaluated PICO element.

**Population:** Population in PICO describes the type of subjects involved in the trial. Critical descriptors for population include important demographic information and any specific shared conditions.

**Intervention:** *Intervention in PICO describes the treatments considered in the trial.* 

**Comparator:** Comparator in PICO describes the alternative treatment to which the intervention is being compared against.

**Outcome:** Outcome in PICO describes the outcome measures used to determine results of the trial. If the primary outcome measures are not be mentioned, then the summary is critically flawed.

The "<*PICOElem*>" tag is replaced with the evaluated PICO element name ("population", "intervention", "comparator", "outcome").

## **H.3 PICO-only LLM Evaluation Prompt**

The prompt used for PICO-only LLM evaluation is almost identical to the one above. The only change is that the first sentence of the previous prompt ("You are given an abstract and a summary.") is changed to "You are given a list of PICO elements from an abstract and a summary."

## **H.4** Prompt for Extracting PICO Elements

The following prompt is used to extract PICO elements from medical text.

Definition of each PICO element:

Population: The types of patients involved in the trial

Intervention: The treatments considered

Comparator: The alternative treatment to which the intervention is being compared to.

Outcome: What is measured. NOT what was found in the study ("result"). For example, if the study finds a drug reduces the duration of headache, the outcome here is just the "duration of headache", not that it reduced it.

Identify the PICO elements in the following

## **H.5** Evidence Inference Full-Text Prompt

In FACTPICO, evaluators analyze if individual evidence inference spans from the abstract are accurately represented in the plain language summary. The LLM evaluation is modeled after this as well, comparing evidence inference spans from the abstract to the full text of the summary. The following is the system prompt used for this evaluation.

You are given a result inference span from an abstract, and you are given a summary. A result inference span corresponds to an inferred result in an experiment. Find the corresponding result inference in the summary and use it to rate the summary between 1 to 4.

The ratings are as follows:

- 1 The result inference is mentioned and described accurately.
- 2 The result inference is mentioned but is described vaguely or is slightly inaccurate.
- 3 The result inference is critically inaccurate.
- 4 The result inference is missing in the model summary.

Please provide only the rating and the rationale for the rating. Provide the rating after stating "Rating:".

The model is then queried as follows:

Result Inference Span: <Span from abstract>

Summary: <Summary>

## **H.6** Evidence Inference Extraction Prompts

We used the following prompt to extract evidence inferences from the plain language summary. These extractions are subsequently compared against evidence inferences from the abstract for the evaluations. The extraction prompt is as follows:

An result inference span corresponds to an inferred result in an experiment.

Identify result inference spans in the following passage. Pull direct quotes from the passage:

Here is the prompt for the evaluation itself:

You are given a result inference span from both an abstract and a summary. A result span corresponds to an inferred result in an experiment. Use the result inference spans from the abstract and the summary to rate the summary between 1 to 4.

The ratings are as follows:

- 1 The result inference is mentioned and described accurately.
- 2 The result inference is mentioned but is described vaguely or is slightly inaccurate.
- *3 The result inference is critically inaccurate.*
- 4 The result inference is missing in the model summary.

Please provide only the rating and the rationale for the rating. Provide the rating after stating "Rating:".

The model is then queried as follows:

Abstract:

< Abstract Evidence Inference Span>

Summary:

< Extracted Evidence Inferences from Summary>

## I Exhaustive Outcomes

The outcome element in RCTs may often be represented through multiple measures, some of which may not be critical for the experiment. The omission of these non-critical outcome measures in plain language summaries usually does not impact its factuality. In FACTPICO we also asked evaluators to determine whether all outcome measures are exhaustively mentioned in the plain language summary as a separate tag exhaustive. This enables us to keep track of when these omissions occur without tying them to the factuality evaluation. The addition of this annotation also enabled better agreement on outcome annotations. This was one of the factors that led to better agreement as discussed in Section C. We also calculated agreement for exhaustive annotations through Randolph's kappa

and report it to be 0.44, which signifies moderate agreement. These collected annotations will also be included in the released data.

#### J Contradictions

	#C	% C
ALPACA	11	8.70
GPT-4	16	10.4
LLAMA-2	36	25.2

Table 12: Total number of contradictions (#C) and percentage of FACTPICO that is a summary with at least one contradiction (%C).

We define *Contradictions* as pieces of information in the plain language summary that meaningfully disagree with the input abstract. In addition to the characteristics described in the main paper, we also ask evaluators to find spans in both the original abstract and the plain language summary that contradict each other and provide a rationale. The results are presented in Table 12, following similar trends to that found for other dimensions in Table 3.

We observed poor agreement overall for this question, largely due to "fuzzy" nature of this task. We define fuzziness as the difficulty to exhaustive distinguish such characteristics within a text. Table 13 shows a fuzzy, contradiction example which demonstrates how certain exaggerations may not be detected by the reader because it may be considered acceptable to use as part of the simplified language style. Annotators with different levels of discernment will likely suffer from low agreement for these types of annotations due to this phenomenon (Trienes et al., 2024). These collected annotations will also be included in the released data.

### **K** Llama-2 Correlation

The negative correlations under the Flip column of Table 14 demonstrate the worser performance of Llama-2 scores that are flipped (post-processed to be 5-x) compared to its original ratings. Out of all LLM evaluators, this occurrence is unique to Llama-2. We attribute this to Llama-2's struggles to remain consistent and faithful to the rating system defined in the prompt. Through analyzing its rationales, we discovered that Llama-2 often ranks summaries it explains as accurate with high numbers, which was opposite to the instructions.

#### **Complete Contradiction**

Source: Compared to OP group, the post operative diclofenac requirement was significantly less in LP group (mean 107.14 mg) and OP group required mean of (682.35 mg).

Simplification [Llama-2]: The study found that laparoscopic pyeloplasty takes longer to do and needs more medicine to help with pain after the surgery, but it also has fewer problems and the person stays in the hospital for less time.

Rationale: The LPA group needed LESS pain medication, not more.

#### Misinterpretation

**Source:** The effects of add-on prolonged-release melatonin (PRM)...were investigated in 80 patients... In this randomized, double-blind, parallel-group study, patients were treated for 2 weeks with placebo and then randomized (1:1) to receive 2 mg of PRM or placebo nightly for 24 weeks, followed by 2 weeks placebo.

Simplification [Llama-2]: The medicine was given to 80 people with AD, some of whom also had trouble sleeping.

Rationale: According to my understanding the 80 people were split 1:1 so that only 40 people were treated with melatonin and 40 received a placebo.

#### Exaggeration (fuzzy)

Source: Although, laparoscopic pyeloplasty has the disadvantages of longer operative time and requires significant skill of intracorporeal knotting but it is here to stay and represents an emerging standard of care.

Simplification [GPT-4]: So, they think the camera surgery is becoming the best choice even if it's a bit more difficult to do.

Rationale: This is not necessarily linked with each other.

Table 13: Examples of contradictions found in plain language summaries within FACTPICO.

This Llama-2 rationale assigns a rating of 4 to a summary that it describes as accurate and sufficiently detailed, which is the opposite of the rating system provided in the prompt.

"The population in the abstract and summary is described accurately and with sufficient detail. The authors mention the type of subjects involved in the trial, including their demographics (age, sex, and duration of HIV infection) and any specific shared conditions (abacavir use). The summary provides a brief overview of the study population, including their baseline characteristics and the results of the intervention. The summary accurately reflects the information provided in the abstract and provides a clear and concise summary of the study population. The authors use descriptive language to highlight the key findings related to the study population, such as the slight improvement in endothelial function after maraviroc intensification. The summary does not contain any inaccuracies or vagueness that could affect the reader's understanding of the study population. Therefore, I would rate the summary as a 4, indicating that the population is described accurately and with sufficient detail."

Additionally, the following rationale demonstrates the hallucinations that Llama-2 makes about the rating scale, as it assumes a ranking order of 1-5 with 5 being the most accurate. Not only is 5 out of the range of the rating options given in the prompt, but also this scale is inconsistent with Llama-2's other evaluation scale of using 4 as its highest rating.

"The population in the abstract and summary is described as "drug users" or "people using illicit opiates," which is a broad and vague term that does not provide critical demographic information. The summary does not provide enough detail about the specific characteristics of the population, such as age, gender, ethnicity, or comorbidities. Additionally, the summary does not indicate whether the population was selected randomly or if there were any exclusion criteria.

To rate the summary, I would give it a score of 4 because it does not provide sufficient information about the population involved in the study. A rating of 5 would indicate that the summary provides detailed and accurate information about the population, including critical demographic information and any specific shared conditions."

	Flip	No Flip
Population	-0.047	0.042
Intervention	-0.063	0.059
Comparator	-0.038	0.036
Outcome	-0.082	0.078
Results	-0.058	0.056
Avg. PICO-R	-0.060	0.055

Table 14: Kendall's Tau correlations between FACT-PICO ratings and the flipped and non-flipped LLAMA-2 ratings.

This analysis explains why non-flipped Llama-2 ratings were more correlated with human evaluation. Llama-2 has a tendency to not follow rating instructions and rate more accurate summaries with higher numbers, defeating the purpose of flipping its ratings in post-processing. Furthermore, Table 5 shows how Llama-2 (non-flipped) is capable of rating plain language summaries similarly to human annotators in Table 3. If we flipped Llama-2 ratings, the average ratings for Llama-2 would have given opposite results. Therefore, we report the non-flipped Llama-2 evaluations throughout this paper to most accurately depict Llama-2's intended ratings.

#### L Rationale Length Analysis

Table 15 shows the average number of tokens for rationales in FACTPICO and those generated by LLMs. As discussed in Section 6, Llama-2 produced the longest rationales. GPT-4, Mistral, and GPT-4 pipelined with PICO-R extraction generated rationales with similar lengths. Among the LLMs, Alpaca produced the shortest rationales. However, overall, expert-written rationales in FACTPICO have the shortest lengths. This is largely because evaluators tend to justify themselves as concisely as possible, especially for easy evaluation instances,

	GPT-4	Llama-2	Alpaca
P	66.6	136.9	33.3
I	64.7	121.4	30.0
C	82.5	140.2	30.1
O	73.8	133.6	28.2
R	63.0	90.8	45.2
	Mistral	Extract	Human
P	66.2	65.1	15.4
I	56.9	55.4	15.0
$\mathbf{C}$	83.4	81.0	12.7
О	66.5	65.4	24.1

Table 15: Average number of tokens for rationales from all systems. Human represents rationales in FACTPICO.

such as when an element is clearly mentioned accurately or clearly missing.

#### **M** LLM Rationale Errors

Table 16 shows several examples of erronous rationales generated by various LLMs. These rationales illustrate several patterns of errors exhibited when these systems generate rationales. The examples from Mistral show the tendency to "forget" to evaluate the summary as described in Section 6. Mistral here either completely ignores the input summary or it passes off text from the abstract as belonging to the summary, as exemplified by the very first example. In addition to this issue, we also observe that Mistral generates illogical and contradictory language in its rationales (ex 3).

The examples from Llama-2 exemplify the length of its rationales as well as its tendency to use irrelevant ideas in its justifications. For example, it may acknowledge factual information from the summary ("While the summary mentions the age range of the patients (adults) and the type of surgery (elective)"), but draws invalid conclusions ("it does not provide additional demographic information").

Rationale examples from GPT-4 here ignores critical factual errors in the summary. For example, the first GPT-4 example shows that the model acknowledges that the summary left critical characteristics missing ("Although the summary doesn't specifically mention 'adult' patients or the type of hospital"), but claims that "this doesn't seem to be crucial information", downplaying its effect during evaluation.

We also included some examples from Alpaca, showing its rationales tend to be very similar to the the rating descriptions provided through the prompt, and exceptions to this tend be nonsensical and illogical text.

## N Full Text of Abstract in Figure 1

Due to space constraints, Figure 1 substituted the full abstract used during evaluation with a link. The full text of this abstract is included below:

NicVAX®, a nicotine vaccine (3'AmNic-rEPA), has been clinically evaluated to determine if higher antibody concentrations are associated with higher smoking abstinence rates and if doses and frequency of administration are associated with increased antibody response. This randomized, double-blinded, placebo-controlled multicenter clinical trial (N=301 smokers) tested 200 and 400 µg doses administered 4 or 5 times over 6 months compared to placebo. 3'AmNic-rEPA recipients with the highest serum anti-nicotine antibody response (top 30% by AUC) were significantly more likely to attain 8 weeks continuous abstinence from weeks 19 through 26 than the placebo recipients (24.6% vs. 12.0%, p=0.024, OR=2.69, 95% CI, 1.14–6.37). The 5 injection 400 µg dose regimen had the greatest antibody response and had significantly higher abstinence rates than placebo. This study demonstrates proof-of-concept that 3'AmNic-rEPA elicits antibodies to nicotine and is associated with higher continuous abstinence rates, justifying its further development as a treatment for nicotine dependence.

Human Rationale: The population of the study is not mentioned clearly at all.

**LLM Rationale [Mistral]:** The population in the abstract and the summary is described accurately. The population includes 16 eyes of 16 patients with corneal neovascularization (NV) secondary to herpetic keratitis, graft rejection, chemical burn, pemphigoid, and recurrent ulcer. The patients were randomly assigned to receive 2.5 mg of bevacizumab or 1 mg of ranibizumab through subconjunctival and intrastromal routes. The researchers followed up with the patients for one month after the injections and calculated the corneal NV areas using Image J software. The results showed that bevacizumab injection resulted in a more effective and stable regression of corneal NV compared to ranibizumab injection.

Human Rationale: It is mentioned that the patients included all (here called vaguely "people") had a kidney transplant, the size of population is missing. LLM Rationale [Mistral]: The population in PICO is mentioned in the abstract and is described accurately. The population includes 90 patients who were prepared for renal transplantation during one year and were randomly classified into two groups: one group with a stent attached to a Foley catheter and the other group with a conventional technique (stent separated from Foley). The population is described in terms of the study design, which is a randomized controlled trial

Human Rationale: Outcome (Low-density lipoprotein cholesterol (LDL-C), PCSK9, and alirocumab levels) is mentioned partly. Missing: PCSK9 LLM Rationale [Mistral]: The outcome measures used in this study are not mentioned in the summary. The summary describes the results of the study, including the reduction in LDL-C levels and the duration of effect, but it does not provide specific information on the outcome measures used to determine these results.

Human Rationale: Comparator (placebo) is not mentioned.

LLM Rationale [Llama-2]: The comparator in the PICO model is the "placebo" condition, which is the control group against which the intervention (THC, CBD, or their combination) is being compared. In the abstract, the comparator is mentioned accurately, but the description is somewhat vague. For example, the authors state that CBD "improves recognition of emotional facial affect," but they do not provide specific details on how this improvement occurs or how it compares to the placebo condition. In the summary, the comparator is mentioned accurately, but the description is still somewhat vague. The authors state that THC alone impairs recognition, while CBD improves it, but they do not provide accurately on the magnitude of these effects or how they compare to the placebo condition. Additionally, the authors mention that the combination of THC and CBD had no effect, but they do not provide any specific data to support this claim. Overall, while the comparator is mentioned accurately in both the abstract and summary, the descriptions provided are not detailed enough to fully understand the results of the study or to compare the interventions being tested.

Human Rationale: Population (200 pat., adults) is mentioned and described accurately.

LLM Rationale [Llama-2]: The population in the summary is described somewhat inaccurately. While the summary mentions the age range of the patients (adults) and the type of surgery (elective), it does not provide additional demographic information such as gender, race, or socioeconomic status. Additionally, the summary does not mention any specific shared conditions that might be relevant to the study. The summary does provide some critical descriptors, such as the types of anesthesia used (separate written consent for anesthesia) and the outcome measures used to assess the efficacy of the informed consent process (adequacy of consent index). However, the summary could benefit from more detailed information about the patient population, such as the number of patients in each group and the duration of follow-up. Overall, while the summary provides some useful information about the study population, it could be improved by providing more detailed demographic information and critical descriptors.

Human Rationale: The population is described almost correctly as "60 healthy adults who are a little overweight". Their diagnose is not "a little" overweight but overweight.

**LLM Rationale [Llama-2]:** The summary provides some information about the population in accordance with PICO, but it is missing critical descriptors in some cases. For example, the summary does not mention the age range of the participants, which is an important demographic characteristic in this study. Additionally, the summary does not provide details about the specific shared conditions among the participants, such as the level of obesity or the presence of metabolic syndrome. While the summary mentions the type of subjects involved in the trial (healthy adults), it does not provide enough detail about the population to accurately rate the summary. Therefore, I will give a rating of 4.

Human Rationale: Population (pat.) is mentioned- yet information about demographics are missing (adults, number of pat.).

**LLM Rationale [GPT-4]:** The population is accurately described in the summary as patients who are going to undergo surgery. This matches the abstract which specifies the population as adult patients undergoing elective surgery at a tertiary care teaching hospital. Although the summary doesn't specifically mention "adult" patients or the type of hospital, this doesn't seem to be crucial information in this context as the aim of the study is more focused on the procedure (anesthesia and consent process) rather than population characteristics. So, the description in the summary is considered accurate.

Human Rationale: The intervention is mentioned vaguely cause it says "one group got the medicine before surgery " - this could refer to any medicine... only in the introduction ketamine is mentioned, yet also it is not made clear that it was an injection of a small doses of Ketamine that was given as intervention.

LLM Rationale [GPT-4]: The abstract and the summary both document the intervention accurately. The intervention mentioned in the abstract is the preemptive use of ketamine before surgery to reduce postoperative pain. The summary accurately describes the same by mentioning that one group got the medicine (ketamine) before surgery, which helped to lessen the pain after the surgery. Both the dosages and the routes of administration have been accurately described in the summary, in line with the abstract's details. The outcomes have also been similarly explained in both. Hence, the summary is rated as 1, indicating that the interventions are mentioned and described accurately.

Human Rationale: Comparator (sham tDCS) is not mentioned and missing.

**LLM Rationale [GPT-4 Extract]:** The summary mentions a form of comparator, represented as "when the other area of the brain called the DLFPC was stimulated in a different way", but it's vaguely related to the original abstract's Comparator of "sham tDCS". The summary fails to accurately explain that the comparison was against a fake (sham) treatment. Therefore, it can be seen as being inaccurately depicted in the summary.

Human Rationale: Intervention (MIC, NPC) is mentioned and described accurately.

LLM Rationale [Alpaca]: The interventions are mentioned in the model summary but described inaccurately.

Human Rationale: The comparator is correctly mentioned as "inpatient vaginal PGE2" IOL.

**LLM Rationale [Alpaca]:** The comparator is described accurately.

Human Rationale: The intervention is correctly described as techniques for removing ureteral stents, linked stent technique.

LLM Rationale [Alpaca]: The linked stent technique has a lower complication rate than the separate stent technique, suggesting its effectiveness and success.

## O Notable Examples

In this section, we show several examples of interest from FACTPICO. We present relevant snippets from the abstract and the appendix and show the rating and rationale given by evaluators.

### **O.1** Obvious Errors

In the following table, we show examples of factual errors that were easily identified by the expert evaluators. They stood out as clear examples of the summary misrepresenting the information in the abstract.

System	Abstract	Summary	Question Type + Rating	Rationale
Alpaca	Patients with inoperable or persistent/recurrent CTEPH (n=261; mean± SD age 59±14 years; 66% women) were randomised to riociguat (up to 2.5 mg three times daily) or placebo.	We compared the effects of riociguat on haemodynamics in patients with inoperable chronic thromboembolic pulmonary hypertension (CTEPH) or persistent/recurrent CTEPH after pulmonary endarterectomy.	Comparator Missing	comparator (placebo) is not mentioned
Alpaca	reatment-naive, non- cirrhotic patients with chronic hepatitis C genotype 6 were treated with pegylated interferon $\alpha$ -2a (180 $\mu$ g/week) and ribavirin (800–1,200 mg, according to weight) for 4 weeks	The study included 152 patients who had an RVR (rapid virological response) and were randomly assigned to either a 24-week or 48-week treatment regimen.	Intervention Missing critical information	The intervention misses the applied medication which is a critical descriptor as the mentioned "treatment regimes" could be anything.
Llama-2	In Group A, 1 µg/kg of dexmedetomidine was intravenously administered for 10 min, followed by the maintenance infusion of dexmedetomidine 0.2 µg/kg/hr after 5 min of intrathecal bupivacaine injection. Patients in Group B received same dose of dexmedetomidine by intravenous administration before 5 min of intrathecal bupivacaine injection.	This study looked at how two different ways of giving a medicine called dexmedetomidine before a spinal anesthesia.	Intervention Described inaccurately	Intervention (timing: dexamedetomidine after spinal anesthesia) is mentioned inaccurately. The original refers to an application before and after spinal anesthesia. But the summary mentions that TWO different ways of giving dexmedetomidine BE-FORE spinal anesthesia were tested.
GPT-4	A total of 84 patients with ureteral calculi 5 to 10 mm in diameter were divided into two groups.	They studied two groups of people who all had this treatment.	Population Missing critical information	The summary only mentiones "They studied two groups of people who all had this treatment". The diagnose is missing, the fact that these are patients and not just people and the size of the group is missing.

Table 17

## **O.2** Subtle Errors

In the following table, we show examples where the error is harder to detect at a first glance. These errors disguise themselves as appearing factual, but a closer reading of the text would reveal factual inconsistencies.

System	Abstract	Summary	Question Type + Rating	Rationale
GPT-4	The sample comprised of forty patients of both sexes diagnosed as acute and transient psychotic disorder who were randomly assigned to high dose (20 mg/day) and low dose (5 mg /day) haloperidol groups with equal number of subjects (n=20) in both groups.	This passage is about a study where they wanted to see if giving bigger or smaller amounts of a medicine(called haloperidol) would work better to help people who had a certain type of brain sickness	Intervention Described vaguely	The intervention is vaguely described as smaller doses from haloperidol, which could be any doses - not even making clear that it is in itself a "low doses".
GPT-4	Seventy-four patients with American Society of Anesthesiologists physical status classification I or II were randomly allocated into two groups.	In this experiment, they divided patients into two groups.	Population Missing critical descriptors	The population (patients) is mentioned. Nevertheless, the specific condition of having an ASA score I / II is not mentioned, which means that the patients are rather healthy and the doctors do not expect so many complications
GPT-4	Two variants near IRS1, rs1522813 and rs2943641, were genotyped in 738 overweight/obese adults (age 60 ± 9 years; BMI 32.7 ± 3.9 kg/m2) randomly assigned to one of four weight-loss diets (a deficit of 750 kcal/day of caloric intake from baseline) varying in macronutrient contents for 2 years	They looked at individuals who were overweight and split them into two groups	Population Described vaguely	The summary only mentions the population as "individuals who were overweight", different facts from the abstract are missing like their genotype, age and BMI.
GPT-4	Group 1 consisted of 14 listeners (number of females=11) who were 41–61 years old [mean=53.18; standard deviation (SD)=5.97].	They tested two groups of adults, one younger and one older, using a special brain test.	Comparator Described vaguely	Comparator (participants aged 41-61 and their performance on the WCJ-III test) is mentioned vaguely as younger people doing a special brain test, which leaves too much space for interpretations.
Llama-2	At day 56, there were more recurrent infections in the chloroquine arm (8.9%, 95% CI 6.0-13.1%) than the dihydroartemisinin-piperaquine arm (2.8%, 95% CI 1.4-5.8%), a difference in cumulative recurrence rate of 6.1% (2-sided 90% CI +2.6 to +9.7%). The log-rank test comparing the survival curves confirmed the superiority of dihydroartemisinin-piperaquine over chloroquine (p = 0.003).	The study found that dihydroartemisinin-piperaquine is better than chloroquine at preventing malaria from coming back after treatment	Outcome Missing critical information	Only recurrence of infection is mentioned as Outcome, other more important like side effects or time till recovery and overall survival are not mentioned.

## O.3 Edge Cases

Here we show examples that evaluators felt were "edge cases." In other words, there is some subjectivity involved in rating these examples according to the FACTPICO framework.

System	Abstract	Summary	Question Type + Rating	Rationale
Llama-2	We compared the haemodynamic effects of riociguat in patients with inoperable chronic thromboembolic pulmonary hypertension (CTEPH) or persistent/recurrent CTEPH	This is a study about a medicine called riociguat. It helps people with a bad heart.	Addition Factual	in first instance CTEPH is a lung disease, that can affect the heart and can cause a "bad heart" as a result.
GPT-4	adult patients, sched- uled for elective cardiac surgery, were divided into three groups: saline, lignocaine (20 mg) and methylprednisolone	They tried mixing the medicine with either saltwater	Comparator Described vaguely	the comparator saline is described vaguely as saltwater.
GPT-4	Changes in the quality of life measured with EQ-5D in the group of men who carried out the PAPP were statistically significant when comparing between groups. However changes in cardiopulmonary function were not as relevant when comparing between groups.	After three months, they found that the men who did the exercise program felt better although there was no change in how they breathed or how their hearts worked. The women did not feel any better.	Outcome Described vaguely	To me the outcome "felt better" as described in the summary seems vague cause it leaves open what a good feeling actually is in this context. The abstract states clearly that it is "quality of life evaluated by a valid questionnaire" - I choose vaguely for this reason.
GPT-4	Intra-articular morphine was as effective as intra- articular triamcinolone for analgesia in patients with osteoarthritis knee. Intra-articular morphine is possibly a better option than intra-articular steroid as it has lesser side effects.	So, if granny or grandpa is having pain because of osteoarthritis, an injection with morphine could help them better.	Addition Non-factual	This is definitely not factual or good advice to give, even though the information within is technically true.

Table 19

#### P Evaluation Interface

FACTPICO evaluations were collected through the Thresh platform (Heineman et al., 2023). This platform allows for the creation of a customizable interface for any text generation task, supporting a variety of question types and text span highlighting. The platform is hosted as a website and the customizable interface is encoded in YAML.

We modify the original Thresh interface to support integration with the Google Drive API. <sup>7</sup> This integration enables annotators to securely retrieve their saved evaluations while having a copy of their evaluations stored in a shared Google drive. This facilitates a straightforward organization of evaluation files in an easily accessible, shared location.

Figure 5 shows both the initial state of the interface and the state after annotations have been completed. The initial state shows almost all questions evaluators would have to answer as "edit annotations." However, span-level annotations have to be created by adding an edit, selecting the type of edit (i.e. Added

<sup>&</sup>lt;sup>7</sup>This version of Thresh is available as a fork of the original project through this link: https://github.com/SebaJoe/thresh

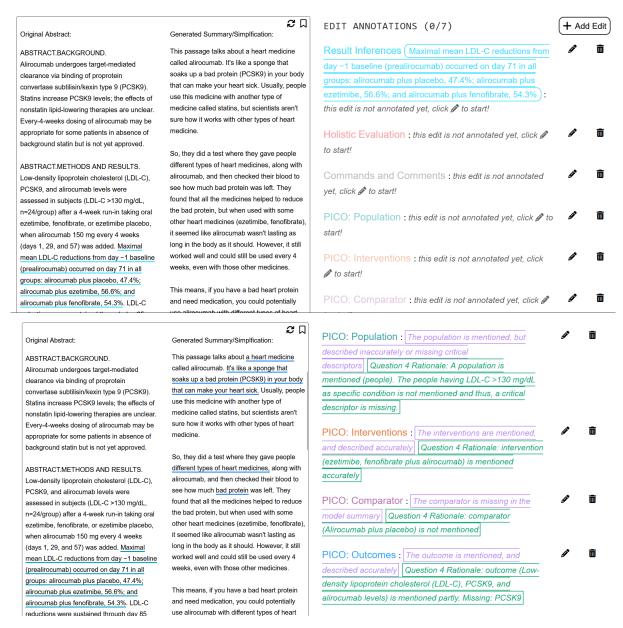


Figure 5: The initial state of the Thresh interface (top) and the state after annotations have been completed (bottom).

Information), and highlighting text. Figure 6 shows how evaluators can select spans in text to evaluate. Evaluators answer questions by clicking the pencil icon next to each edit. Figure 7 shows how questions are displayed for evaluators to answer.

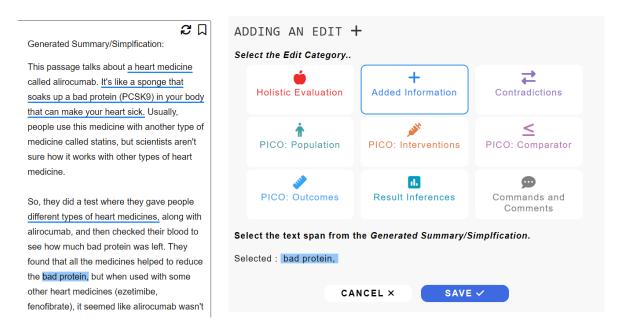


Figure 6: Annotating an added information span in a plain language summary.

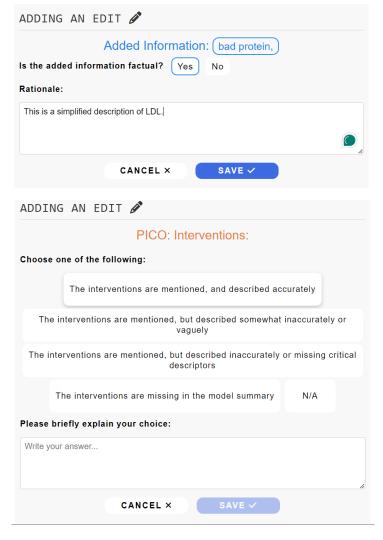


Figure 7: Interface for answering questions regarding added information (top) and PICO interventions (bottom).