

# FrontierScience Bench: Evaluating AI Research Capabilities in LLMs

Matthew Li\*, Santiago Torres-Garcia\*, Shayan Halder, Phani Kuppa,  
Vasu Sharma, Sean O’Brien, Kevin Zhu, Sunishchal Dev

Algoverse AI Research

santiago\_122303@berkeley.edu, kevin@algoverse.us, sean@algoverse.us

## Abstract

Large language models (LLMs) have shown remarkable capabilities across various tasks, yet their potential to reason about and construct scientific methodologies remains under explored. This work introduces a novel benchmark evaluating LLMs’ capacity to predict methodological details in AI research papers. We construct a dataset of 88 papers with redacted methodology sections and zero-shot prompt several state-of-the-art LLMs to generate methodology predictions. Our evaluation framework then employs a LLM-as-judge system with multiple LLM judges, majority voting, and self-omission techniques to minimize biases. We validate our LLM judge scores against human judgments. We then briefly analyze the judging results of our zero-shot prediction pipeline, suggesting that even state-of-the-art LLMs struggle with the task of methodology generation without more advanced techniques. This benchmark lays the groundwork for future research into evaluating LLMs’ potential for aiding in AI research. Our benchmark code and dataset are open-sourced at <https://github.com/Swadian/FrontierScience-Bench>

## 1 Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities in tasks such as summarization, translation, and text generation (Fabri et al., 2020; Zhu et al., 2023; Ni and Li, 2024). While LLMs excel at identifying patterns from their vast training data, their ability to engage in the structured reasoning required for scientific discovery remains underexplored. This is crucial because in many scientific fields, progress relies not just on understanding existing knowledge but on devising novel methodologies to transform abstract goals into rigorous experimental designs. To

test LLMs’ capacity for scientific reasoning, we introduce a benchmark that evaluates their ability to infer plausible methodology sections from redacted AI research papers, where explicit methodological details have been removed. Our study provides a high-quality dataset of 88 manually redacted papers, a zero-shot prediction pipeline, and a rigorous LLM-as-a-Judge evaluation framework to assess the quality of predictions.

## 2 Related Works

The application of LLMs as assistants in the scientific method is a promising research area. Existing benchmarks for LLM evaluation, such as factual knowledge retrieval, summarization, or question answering, focus mainly on surface-level tasks (Guo et al., 2023; McIntosh et al., 2024; Porcu and Havlínová, 2024). These benchmarks test recall and synthesis but rarely assess whether LLMs can reason creatively and simulate problem-solving processes essential for scientific discovery.

Luo et al. (2024b) demonstrated that LLMs have surpassed experts in predicting neuroscience experiment outcomes through their benchmark, BrainBench. While outcome prediction is valuable, it reveals little about an LLM’s ability to plan and reason through methodological steps. This abstract reasoning process is crucial for aiding researchers in accelerating scientific discovery. Sun et al. (2024b) developed the SciEval benchmark to examine LLM reasoning in biology, chemistry, and physics. However, its multiple-choice format primarily assesses textbook knowledge and calculation skills rather than extensive reasoning. While useful for evaluating general understanding, these questions do not push LLMs to plan and organize methodological steps—a critical skill when considering the background and context of a research study.

\*Equal contribution.

Acknowledgment to Jason Wu for his help with dataset refinement and prompting.

Gu et al. (2025a) systematically examines the challenges of using LLMs for evaluation tasks, emphasizing reliability, bias mitigation, and resistance to manipulation. Their survey addresses critiques of LLM-as-a-Judge systems, including biases like length bias and self-enhancement bias, and discusses mitigation approaches such as multiple LLM judges, structured outputs, voting-based consensus mechanisms, and fine-tuning. While their study covers broad applications, our work focuses specifically on methodology generation and assessing the accuracy of these methodologies.

Weng et al. (2025) uses LLM agents to simulate the full research–review–refinement workflow, while Starace et al. (2025) evaluates AI agents’ ability to replicate 20 ICML 2024 papers via hierarchical rubrics and LLM-based judging—finding only about 21 % average replication. In contrast, our benchmark isolates zero-shot methodology prediction for AI/ML papers and employs a multi-judge LLM framework with majority voting and self-omission to reduce bias. This narrow focus complements broader end-to-end and replication studies by targeting the core methodological reasoning essential for novel AI research.

### 3 Methods

In this section we explain the details behind the construction of our benchmark. As a reminder, our goal is to evaluate the extent to which LLMs can reconstruct novel research methodologies based on redacted AI research papers.

The benchmark consists of three stages: dataset curation, prediction, and evaluation. We first created a dataset of 88 research papers, redacted to remove their methodology sections, results, and any references to methodology found in other sections. We then used this new dataset and zero-shot prompted multiple state-of-the-art LLM’s including GPT-4o, o3-mini, Claude 3.5 Sonnet, and Gemini 1.5 Pro to generate methodologies. Finally, we used a LLM-as-a-Judge framework to evaluate how close the predicted methodologies were to the original ones. Detailed prompts, examples, and additional figures are provided in the appendices (D, E, F).

#### 3.1 Original Paper Curation

To build our dataset, we curated 88 research papers from top conferences such as EMNLP, ACL, IEEE, AAAI, ICML, and CHI using Semantic Scholar.

All published in 2024, after the knowledge cutoff of our initial predictor, GPT-4o, (October 2023) (OpenAI, 2023). This timing minimizes training data contamination, though our benchmark will need future updates as newer LLMs with more recent cutoffs (e.g., Claude 3.5 Sonnet and Gemini 1.5 Pro, with cutoffs in April and May 2024, respectively (Anthropic, 2024; Cloud, 2024) are released, which we discuss in Section 6.

We automated paper collection using Semantic Scholar’s API, filtering for 2024 conference papers with the keywords “machine learning” and “large language models.” We focused on this domain both because our expertise facilitates better redaction and validation, and due to our interest in recursive self-improvement in AI systems—a capability linked to potential intelligence explosions (Barrett and Baum, 2016). Papers were limited to 15 pages to comply with context length restrictions, and we parsed PDFs using PyMuPDF (PyMuPDF Developers, 2024), excluding those with extensive mathematical notation due to issues parsing  $\LaTeX$ . See Appendix E for an example.

Once the papers were collected, we manually removed irrelevant content (e.g., figures, headers, footers, references, appendices) while retaining table and figure captions. We then prompted GPT-4o with a `rewrite_paper_prompt` and enabled Structured Outputs (OpenAI, 2023) to convert each paper into a JSON format with keys for the abstract, introduction, related works, and methodology—yielding our *trimmed paper*. We later extracted the original methodologies from these trimmed papers for evaluation. A visualization of our collection process is provided in Appendix F.

#### 3.2 Redaction Pipeline

After selecting and processing papers, we passed them through our two-stage redaction pipeline to remove the author’s methodologies. We provide a visualization of this pipeline in Appendix F. This first involved taking the trimmed paper in JSON format and removing the methodologies field (denoted as *filtered paper*). Then, to deal with minor instances of the methodologies that were spread throughout the remaining sections, we had two layers of manual removal by the authors of this paper. The first step had each author remove revealing information from their assigned range of papers. The second step then had a single author review all of these annotations, alongside the trimmed paper and make changes where necessary. This helped

to ensure high-fidelity redaction and reduced variability from having many human annotators. To aid the process of manual redaction, we used DiffChecker to visualize differences between original and redacted papers.

A limitation of our manual approach is the burden of human redaction, despite our efforts to increase automation. We believe automation is crucial for creating new benchmark datasets for future models with later knowledge cutoffs. We describe in-detail three approaches toward automation (and why they failed) in Appendix A.

### 3.3 Prediction

After the redacted papers were constructed, we proceeded with our prediction phase, where a predictor LLM was prompted with the redacted paper to generate a methodology. Our method for this phase was a simple zero-shot baseline that we entirely credit to Si et al. (2024). Namely, we first prompted the predictor LLM (with structured outputs enabled) (OpenAI, 2023; Cloud, 2024; Anthropic, 2024) with the redacted paper and a `outline_prompt`, which generated a `proposed_method` that highlights all the necessary steps of the method and also a `experimental_plan`, which further elaborates on the steps in the `proposed_method`, covering more details like specific models and datasets. Then, this outline is used as input once again for the predictor LLM, this time using a `writing_prompt`, to generate a cohesive written methodology, mimicking the style that might be found in a research paper.

### 3.4 Evaluation

After collecting the methodology predictions, the next step was to evaluate the quality of these predictions. We evaluated the methodology predictions using an open-form LLM-as-a-Judge approach, which is more scalable compared to human-only evaluations and better aligned with human judgment compared to closed-form approaches like multiple-choice. For completeness, we provide an overview of our attempts at MCQ in Appendix C. While there have been several critiques of LLM-as-a-Judge in the past as we mentioned in Section 2, we argue that our LLM judging pipeline is more aligned with humans and mitigates pitfalls such as bias in the judging process. We visualize our pipeline in Appendix F.

Our evaluation pipeline uses multiple LLMs and majority voting to reduce variability in LLM-as-

a-Judge systems (Gu et al., 2025b). A jury of five LLMs conducts majority voting across five runs, with final scores averaged. To mitigate self-preference bias (Gu et al., 2025b), we exclude predictor LLMs from judging. Judges use chain-of-thought reasoning (Wei et al., 2022) and a detailed rubric to score predicted vs. original methodologies from 1-10. See Appendix D for details.

## 4 Experiments and Results

Our framework uses GPT-4o, o3-mini, Claude 3.5 Sonnet, and Gemini 1.5 Pro as predictors, while judges include GPT-4o, o3-mini, Claude 3.5 Sonnet, and Grok 2. Models were chosen for their strong instruction-following capabilities. Grok 2 was excluded from prediction due to inconsistent outputs and uncertain cutoff; Gemini 1.5 Pro from judging due to bias. Deepseek R1 and Llama 3.1 405B Instruct failed initial trials due to JSON formatting issues. We provide examples and extra figures in Appendix E and Appendix F.

### 4.1 Results Analysis

Figure 1 shows the aggregated score distributions for each predictor LLM, illustrating overall performance trends. All models average between 3 and 3.5. Due to the knowledge cutoff, we refrain from directly comparing the predictor LLMs. Nonetheless, the right-skewed distributions and low averages suggest that state-of-the-art LLMs struggle to consistently produce high-quality methodologies when prompted zero-shot. Although they capture some fundamental steps, they lack the fine-grained, domain-specific details needed for rigor. These findings highlight the need for further innovation to guide LLMs toward generating robust scientific methodologies. Specific examples of these predictions are provided in Appendix E. Statistical significance and power analysis is discussed in Appendix B.

### 4.2 Evidence Toward Reliable Evaluation

To validate our evaluation process, we create scatter plots comparing predicted methodology length and final judge scores. We visualize this in Appendix F. Visually, we observe no correlation between the two variables, suggesting that length bias is mitigated and reinforcing our method’s reliability.

We also compare the jury’s aggregated scores with those of a human evaluator across a validation set of 10 predicted methodologies. We calculate

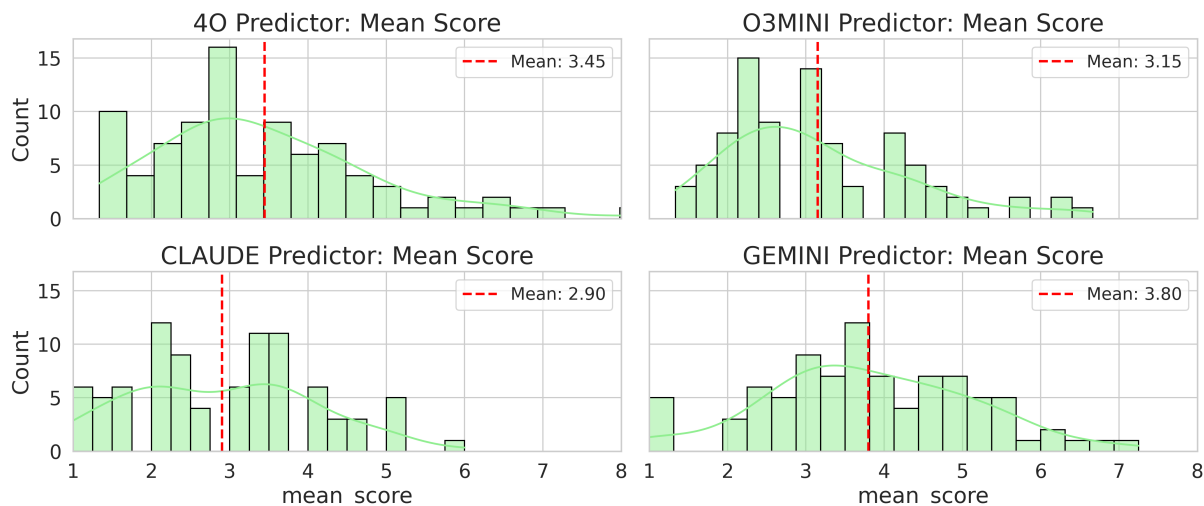


Figure 1: Comparison of the aggregate score distributions across various predictor LLMs. Each histogram represents the frequency of final scores (x-axis) assigned to the generated methodologies across the dataset of research papers. The red vertical line in each figure denotes the mean final scores.

the absolute error and mean absolute error, where a lower score indicates closer alignment with human judgment. As shown in Table 1, the absolute error was minimal, with nearly all methodologies differing by less than 1. The mean absolute error of 0.725 confirms our evaluation closely aligns with human judgment.

## 5 Conclusion

In this work, we introduced a novel benchmark to evaluate the ability of large language models to generate detailed scientific methodologies from redacted research papers. Our dataset of 88 redacted papers, combined with a zero-shot prediction pipeline and a robust LLM-as-a-Judge evaluation framework, demonstrates that while state-of-the-art LLMs can capture broad methodological outlines, they consistently struggle to deliver the domain-specific precision required for high-quality scientific methods. The close alignment with human evaluations (mean absolute error of 0.725) reinforce the reliability of our evaluation process.

Future work should expand the dataset to more scientific disciplines and better handle mathematical notation. Automating paper redaction will help publish updated versions of the dataset for newer LLMs. AI agents with code execution tools will allow models to validate predicted methodologies before scoring. Finally, including more LLM judges with expanded human validation will increase reliability in scoring.

ID	Agg. Score	Human Score	Abs. Error
000_v1	3.50	4	0.50
004_v1	3.25	4	0.75
005_v1	5.50	4	1.50
009_v1	5.00	4	1.00
011_v1	5.25	4	1.25
000_v2	8.25	8	0.25
004_v2	7.50	8	0.50
005_v2	8.25	8	0.25
009_v2	8.00	8	0.00
011_v2	8.25	8	0.25

Table 1: Aggregated LLM jury scores, human score, and absolute errors for proposed methodologies in validation set. Mean Abs. Error (MAE) = 0.625

## 6 Limitations

### 6.1 Paper Collection

Our dataset consist of 88 research papers exclusively in the AI/ML domain with 15 pages or less, excluding references and appendices. However, our restrictions on these papers bring about some concerns. Although this domain was chosen because of the research team’s familiarity, it narrows the scope of our evaluation, making it unclear whether our findings would generalize to other disciplines, such as medicine or social sciences, where experimental frameworks may differ significantly.

We imposed this page limit due to context window limits and to constrain the manual redaction process. This may exclude more rigorous method-

ologies in longer research papers that are more challenging to predict. Additionally, papers with extensive mathematical notation were excluded due to PyMuPDF parsing failures. This limits our benchmark’s applicability to theoretical or math-heavy research and may inflate model performance. Finally, while all 88 papers were published after GPT-4o’s October 2023 cutoff to minimize training data overlap, we later included models with later cutoffs (e.g., Claude 3.5 Sonnet, Gemini 1.5 Pro). This introduces a risk of training data contamination in those models.

## 6.2 Redaction Pipeline

Our redaction pipeline, which removes explicit methodology details from research papers to prevent LLMs from memorizing and regurgitating author contributions, introduces several limitations. A key challenge is that methodologies are often interwoven into sections like the introduction and related works, creating the risk of unintentional leakage. Despite implementing a two-pass manual review process, subtle methodological hints may remain, undermining the integrity of our dataset. Additionally, manual redaction is also time-consuming and not scalable for future benchmark updates. Despite our efforts toward an automated redaction process, which is briefly described in Section 3 and further elaborated in Appendix A, we fail to find an approach that is both accurate and consistent.

## 6.3 Methodology Evaluation

A key concern with our approach is that it focuses solely on comparing predicted methodologies to the original methodologies, without assessing their practical outcomes. It is possible that an LLM-generated methodology, while different in design, could yield results comparable to or even superior to those of the original research. By relying solely on textual similarity to the ground truth, our evaluation may overlook valid, creative approaches that demonstrate comparable scientific reasoning. Additionally, our evaluation process does not incorporate experimental validation, which would assess whether the predicted methodologies could reproduce the outcomes reported in the original research. This step is especially important given the tendency of LLMs to hallucinate, introducing the risk of false or unrealizable methodologies. As a result, our current evaluation may both understate the reasoning ability of LLMs and fail to catch critical

flaws in their output.

## 6.4 LLM-as-a-Judge Systems

While our LLM-as-a-Judge framework employed several techniques to mitigate common biases and enhance evaluation consistency, it is not without limitations. One limitation is the lack of sufficient human oversight in our evaluation process. While we developed a validation set and tested against a singular human judge to verify the alignment of LLM judges with human evaluators, this may not be sufficient to draw conclusions about LLM-judge alignment with human evaluators. Although LLM-based evaluation offers scalability, significant improvements to LLM-as-a-Judge systems still fail to achieve consistent alignment with human evaluators. (Gu et al., 2025b; Chen et al., 2024a). However, it should be noted that human-alignment should not always be considered a perfect metric, as human evaluators are also prone to bias and variability in open-ended tasks such as text-evaluation.

Furthermore, while majority voting (majority@5) was selected for its superior ability compared to other techniques outlined in Gu et al. (2025b) such as taking the mean score(- mean@5), and taking the best score(- best@5), conducting 5 trials per LLM judge may not be sufficient to get a representative consensus. Lastly, a limitation of our LLM-as-a-Judge framework is the relatively small number of judges employed. Many of our evaluations were conducted with only 3 judge LLMs, which may limit the robustness of our aggregated scores.

## 6.5 Ethics

This benchmark is designed as a diagnostic tool to assess the reasoning capabilities of LLMs and does not aim to automate scientific authorship or replace human researchers. This task isolates the problem of reconstructing plausible methodological reasoning from surrounding content for evaluation purposes. We acknowledge the risks of misuse, such as treating generated methods as ready for use in real AI research. To reduce this risk, we frame our task as a diagnostic benchmark, not a writing tool, and do not recommend using these systems in scientific work without safeguards.

## References

- Anthropic. 2024. How up-to-date is claude’s training data?
- Anthony Michael Barrett and Seth D. Baum. 2016. A model of pathways to artificial superintelligence catastrophe for risk and decision analysis. *Journal of Experimental & Theoretical Artificial Intelligence*, 29:397 – 414.
- Francisco El nio Bezerra et al. 2024. Impacts of feature selection on predicting machine failures by machine learning algorithms. *Applied Sciences*.
- Paul Calle et al. 2024. Towards ai-driven healthcare: Systematic optimization, linguistic analysis, and clinicians’ evaluation of large language models for smoking cessation interventions. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. ACM.
- Daniele Rege Cambrin et al. 2024. Beyond accuracy optimization: Computer vision losses for large language model fine-tuning. *Preprint*, arXiv:2409.13641.
- Yihan Cao et al. 2024. Tarot: A hierarchical framework with multitask co-pretraining on semi-structured data towards effective person-job fit. *Preprint*, arXiv:2401.07525.
- Guiming Hardy Chen, Shunian Chen, Ziche Liu, Feng Jiang, and Benyou Wang. 2024a. Humans or llms as the judge? a study on judgement bias.
- Liuqing Chen et al. 2024b. Chatscratch: An ai-augmented system toward autonomous visual programming learning for children aged 6-12. *Preprint*, arXiv:2402.04975.
- Xiang Chen et al. 2024c. Unified hallucination detection for multimodal large language models. *Preprint*, arXiv:2402.03190.
- Pranjal A. Chitale et al. 2024. An empirical study of in-context learning in llms for machine translation. *Preprint*, arXiv:2401.12097.
- Ye-eun Cho and Seong mook Kim. 2024. Pragmatic inference of scalar implicature by llms. *Preprint*, arXiv:2408.06673.
- Bryan R Christ et al. 2024. Mathwell: Generating educational math word problems using teacher annotations. *Preprint*, arXiv:2402.15861.
- Google Cloud. 2024. Gemini 1.5 pro - vertex ai documentation.
- Clayton Cohn, Nicole Hutchins, and Tuan Le. 2024a. A chain-of-thought prompting approach with llms for evaluating students’ formative assessment responses in science. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38:23182–23190.
- Clayton Cohn et al. 2024b. A chain-of-thought prompting approach with llms for evaluating students’ formative assessment responses in science. *Preprint*, arXiv:2403.14565.
- Tamara Czinczoll et al. 2024. Nextlevelbert: Masked language modeling with higher-level representations for long documents. *Preprint*, arXiv:2402.17682.
- Luis Miguel Vieira da Silva et al. 2024. Toward a method to generate capability ontologies from natural language descriptions. *Preprint*, arXiv:2406.07962.
- Naihao Deng et al. 2024. Tables as texts or images: Evaluating the table reasoning ability of llms and mllms. *Preprint*, arXiv:2402.12424.
- Akash Dutta and Ali Jannesari. 2024. Miren-coder: Multi-modal ir-based pretrained embeddings for performance optimizations. *Preprint*, arXiv:2407.02238.
- Jessica Echterhoff et al. 2024. Cognitive bias in decision-making with llms. *Preprint*, arXiv:2403.00811.
- A. R. Fabbri, Wojciech Kryscinski, Bryan McCann, Richard Socher, and Dragomir R. Radev. 2020. Summeval: Re-evaluating summarization evaluation. *Transactions of the Association for Computational Linguistics*, 9:391–409.
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, Saizhuo Wang, Kun Zhang, Yuanzhuo Wang, Wen Gao, Lionel Ni, and Jian Guo. 2025a. A survey on llm-as-a-judge. *arXiv preprint*, 2411.15594.
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, Saizhuo Wang, Kun Zhang, Yuanzhuo Wang, Wen Gao, Lionel Ni, and Jian Guo. 2025b. A survey on llm-as-a-judge. *Preprint*, arXiv:2411.15594.
- Zishan Guo, Renren Jin, Chuang Liu, Yufei Huang, Dan Shi, Supryadi, Linhao Yu, Yan Liu, Jiaxuan Li, Bojian Xiong, and Deyi Xiong. 2023. Evaluating large language models: A comprehensive survey. *ArXiv*, abs/2310.19736.
- Abhay Gupta et al. 2024. Aavenue: Detecting llm biases on nlu tasks in aave via a novel benchmark. *Preprint*, arXiv:2408.14845.
- Shabnam Hassani et al. 2024. Rethinking legal compliance automation: Opportunities with large language models. *Preprint*, arXiv:2404.14356.
- Chi-Yang Hsu et al. 2024. Thought graph: Generating thought process for biological reasoning. *Preprint*, arXiv:2403.07144.
- Daniel P. Jeong et al. 2024. Medical adaptation of large language and vision-language models: Are we making progress? *Preprint*, arXiv:2411.04118.

- Zhen Jia et al. 2024. [Tiq: A benchmark for temporal question answering with implicit time constraints](#). In *Companion Proceedings of the ACM Web Conference 2024*. ACM.
- Mollie Jordan et al. 2024. [Need a programming exercise generated in your native language? chatgpt's got your back: Automatic generation of non-english programming exercises using openai gpt-3.5](#). In *Proceedings of the 55th ACM Technical Symposium on Computer Science Education v. 1*. ACM.
- Athanasios Karapantelakis et al. 2024. [Using large language models to understand telecom standards](#). *Preprint*, arXiv:2404.02929.
- Minsang Kim et al. 2024a. [Qpaug: Question and passage augmentation for open-domain question answering of llms](#). *Preprint*, arXiv:2406.14277.
- Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoon Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2024b. [Prometheus: Inducing fine-grained evaluation capability in language models](#). *Preprint*, arXiv:2310.08491.
- Taewan Kim et al. 2024c. [Diarymate: Understanding user perceptions and experience in human-ai collaboration for personal journaling](#). In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. ACM.
- Brendan King and Jeffrey Flanigan. 2024. [Unsupervised end-to-end task-oriented dialogue with llms: The power of the noisy channel](#). *Preprint*, arXiv:2404.15219.
- Klaus Krippendorff. 2011. [Computing krippendorff's alpha-reliability](#).
- Shreyas Kumar et al. 2024. [What the phish! effects of ai on phishing attacks and defense](#). *International Conference on AI Research*.
- Daniil Laptev, Nikita Balagansky, Yaroslav Aksenov, and Daniil Gavrilov. 2025. [Analyze feature flow to enhance interpretation and steering in language models](#). *Preprint*, arXiv:2502.03032.
- Bruce W. Lee and JaeHyuk Lim. 2024. [Language models don't learn the physical manifestation of language](#). *Preprint*, arXiv:2402.11349.
- Joshua Lee et al. 2024a. [Pragmatic metacognitive prompting improves llm performance on sarcasm detection](#). *Preprint*, arXiv:2412.04509.
- Yoonjoo Lee et al. 2024b. [Paperweaver: Enriching topical paper alerts by contextualizing recommended papers with user-collected papers](#). *Preprint*, arXiv:2403.02939.
- Florian Leiser et al. 2024. [Hill: A hallucination identifier for large language models](#). *Preprint*, arXiv:2403.06710.
- Bohao Li, Yuying Ge, et al. 2023. [Seed-bench-2: Benchmarking multimodal large language models](#). *Preprint*, arXiv:2311.17092.
- Fan Li et al. 2024a. [Virtual co-pilot: Multimodal large language model-enabled quick-access procedures for single pilot operations](#). *Preprint*, arXiv:2403.16645.
- Weijie Li, Zhentao Gu, et al. 2024b. [Exploring large language models text style transfer capabilities](#). *IOS Press*.
- Zhigen Li, Yanmeng Wang, et al. 2024c. [Learning to adapt to low-resource paraphrase generation](#). *Preprint*, arXiv:2412.17111.
- Zhuochun Li, Bo Xie, et al. 2024d. [Effects of different prompts on the quality of gpt-4 responses to dementia care questions](#). *Preprint*, arXiv:2404.08674.
- Haokun Lin et al. 2024. [Mope-clip: Structured pruning for efficient vision-language models with module-wise pruning error metric](#). *Preprint*, arXiv:2403.07839.
- Hanchao Liu et al. 2024a. [Programmable motion generation for open-set motion control tasks](#). *Preprint*, arXiv:2405.19283.
- Zhengyuan Liu, Stella Xin Yin, et al. 2024b. [Scaffolding language learning via multi-modal tutoring systems with pedagogical instructions](#). *Preprint*, arXiv:2404.03429.
- Zhicheng Liu, Ali Braytee, et al. 2024c. [Ensemble pretrained models for multimodal sentiment analysis using textual and video data fusion](#). In *Companion Proceedings of the ACM Web Conference 2024*. ACM.
- Do Xuan Long et al. 2024. [Multi-expert prompting improves reliability, safety, and usefulness of large language models](#). *Preprint*, arXiv:2411.00492.
- Keming Lu et al. 2024. [Large language models are superpositions of all characters: Attaining arbitrary role-play via self-alignment](#). *Preprint*, arXiv:2401.12474.
- Chuwei Luo et al. 2024a. [Layoutllm: Layout instruction tuning with large language models for document understanding](#). *Preprint*, arXiv:2404.05225.
- Xiaoliang Luo, Akilles Rechart, Guangzhi Sun, Kevin K. Nejad, Felipe Yáñez, Bati Yilmaz, Kangjoo Lee, Alexandra O. Cohen, Valentina Borghesani, Anton Pashkov, Daniele Marinazzo, Jonathan Nicholas, Alessandro Salatiello, Ilia Sucholutsky, Pasquale Minervini, Sepehr Razavi, Roberta Rocca, Elkhani Yusifov, Tereza Okalova, Nianlong Gu, Martin Ferrianc, Mikail Khona, Kaustubh R. Patil, Pui-Shee Lee, Rui Mata, Nicholas E. Myers, Jennifer K. Bizley, Sebastian Musslick, Isil Poyraz Bilgin, Guiomar Niso, Justin M. Ales, Michael Gaebler, N. Apurva Ratan Murty, Leyla Loued-Khenissi, Anna Behler, Chloe M. Hall, Jessica Dafflon, Sherry Dongqi Bao,

- and Bradley C. Love. 2024b. Large language models surpass human experts in predicting neuroscience results.
- Zilin Ma et al. 2024. Evaluating the experience of lgbtq+ people using large language model based chatbots for mental health support. *Preprint*, arXiv:2402.09260.
- Rangaswamy Madugundu et al. 2024. Optimal timing of carrot crop monitoring and yield assessment using sentinel-2 images: A machine-learning approach. *Applied Sciences*.
- Dimitrios Michael Manias et al. 2024. Towards intent-based network management: Large language models for intent extraction in 5g core networks. *Preprint*, arXiv:2403.02238.
- Timothy R. McIntosh, Teo Susnjak, Tong Liu, Paul Waters, and Malka N. Halgamuge. 2024. Inadequacies of large language model benchmarks in the era of generative artificial intelligence. *ArXiv*, abs/2402.09880.
- Luise Metzger et al. 2024. Empowering calibrated (dis)trust in conversational agents: A user study on the persuasive power of limitation disclaimers vs. authoritative style. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. ACM.
- Jamshid Mozafari et al. 2024. Exploring hint generation approaches in open-domain question answering. *Preprint*, arXiv:2409.16096.
- Vy Nguyen and Chau Pham. 2024. Leveraging large language models for suicide detection on social media with limited labels. *Preprint*, arXiv:2410.04501.
- Xuanfan Ni and Piji Li. 2024. A systematic evaluation of large language models for natural. *ArXiv*, abs/2405.10251.
- Daniel Nichols et al. 2024. Can large language models write parallel code? *Preprint*, arXiv:2401.12554.
- OpenAI. 2023. Openai documentation.
- Liyizhe Peng et al. 2023. Customising general large language models for specialised emotion recognition tasks. *Preprint*, arXiv:2310.14225.
- Nisha Pillai et al. 2024. Endtoendml: An open-source end-to-end pipeline for machine learning applications. *Preprint*, arXiv:2403.18203.
- Valentina Porcu and Aneta Havlínová. 2024. Breaking down the metrics: A comparative analysis of llm benchmarks. *International Journal of Science and Research Archive*.
- PyMuPDF Developers. 2024. Pymupdf documentation.
- Federico Ranaldi et al. 2024. Investigating the impact of data contamination of large language models in text-to-sql translation. *Preprint*, arXiv:2402.08100.
- Zafaryab Rasool et al. 2024. Llms for test input generation for semantic caches. *Preprint*, arXiv:2401.08138.
- Rajat Rawat et al. 2024. Diversitymedqa: Assessing demographic biases in medical diagnosis using large language models. *Preprint*, arXiv:2409.01497.
- Matthew Renze and Erhan Guven. 2024. The effect of sampling temperature on problem solving in large language models. *Preprint*, arXiv:2402.05201.
- Paul Röttger et al. 2024. Political compass or spinning arrow? towards more meaningful evaluations for values and opinions in large language models. *Preprint*, arXiv:2402.16786.
- Joni Salminen et al. 2024. Deus ex machina and personas from large language models: Investigating the composition of ai-generated persona descriptions. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. ACM.
- Soma Sato et al. 2024. Improving sentence embeddings with automatic generation of training data using few-shot examples. *Preprint*, arXiv:2402.15132.
- Yash Saxena et al. 2024. Evaluating consistency and reasoning capabilities of large language models. *Preprint*, arXiv:2404.16478.
- Nikhil Sharma et al. 2024a. Generative echo chamber? effects of llm-powered search systems on diverse information seeking. *Preprint*, arXiv:2402.05880.
- Roshan Sharma et al. 2024b. Speech vs. transcript: Does it matter for human annotators in speech summarization? *Preprint*, arXiv:2408.07277.
- Ying Sheng et al. 2024. Measuring an llm's proficiency at using apis: A query generation strategy. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. ACM.
- Wenqi Shi et al. 2024a. Ehragent: Code empowers large language models for few-shot complex tabular reasoning on electronic health records. *Preprint*, arXiv:2401.07128.
- Yuchen Shi et al. 2024b. Agentre: An agent-based framework for navigating complex information landscapes in relation extraction. *Preprint*, arXiv:2409.01854.
- Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. 2024. Can llms generate novel research ideas? a large-scale human study with 100+ nlp researchers. *Preprint*, arXiv:2409.04109.
- Aryan Singhal, Veronica Shao, Gary Sun, Ryan Ding, Jonathan Lu, and Kevin Zhu. 2024a. "a comparative study of translation bias and accuracy in multilingual large language models for cross-language claim verification".



- Aryan Singhal et al. 2024b. [A comparative study of translation bias and accuracy in multilingual large language models for cross-language claim verification](#). *Preprint*, arXiv:2410.10303.
- G. Starace, O. Jaffe, D. Sherburn, J. Aung, J. S. Chan, L. Maksin, R. Dias, E. Mays, B. Kinsella, W. Thompson, J. Heidecke, A. Glaese, and T. Patwardhan. 2025. [PaperBench: Evaluating AI's Ability to Replicate AI Research](#). arXiv preprint arXiv:2504.01848. Version posted April 7, 2025.
- Jiaxing Sun et al. 2024a. [Benchmarking chinese commonsense reasoning of llms: From chinese-specifics to reasoning-memorization correlations](#). *Preprint*, arXiv:2403.14112.
- Liangtai Sun, Yang Han, Zihan Zhao, Da Ma, Zhennan Shen, Baocai Chen, Lu Chen, and Kai Yu. 2024b. [Scieval: A multi-level large language model evaluation benchmark for scientific research](#). *Preprint*, arXiv:2308.13149.
- Wannapon Suraworachet et al. 2024. [Predicting challenge moments from students' discourse: A comparison of gpt-4 to two traditional natural language processing approaches](#). *Preprint*, arXiv:2401.01692.
- Annalisa Szymanski et al. 2024. [Integrating expertise in llms: Crafting a customized nutrition assistant with refined template instructions](#). In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. ACM.
- Mahsa Tavasoli, Euihark Lee, Yashar Mousavi, Hananeh Barahouei Pasandi, and Afef Fekih. 2024. [Wipe: A novel web-based intelligent packaging evaluation via machine learning and association mining](#). *IEEE Access*, 12:45936–45947.
- Bekzat Tilekbay et al. 2024. [Expressedit: Video editing with natural language and sketching](#). *Preprint*, arXiv:2403.17693.
- Eugenia I. Toki et al. 2024. [Machine learning for predicting neurodevelopmental disorders in children](#). *Applied Sciences*.
- Mesut Ulu et al. 2024. [Prediction of traffic incident locations with a geohash-based model using machine learning algorithms](#). *Applied Sciences*.
- Bin Wang et al. 2024a. [Resilience of large language models for noisy instructions](#). *Preprint*, arXiv:2404.09754.
- Che Wang et al. 2024b. [Contracttinker: Llm-empowered vulnerability repair for real-world smart contracts](#). *Preprint*, arXiv:2409.09661.
- Keshara Weerasinghe et al. 2024. [Real-time multimodal cognitive assistant for emergency medical services](#). *Preprint*, arXiv:2403.06734.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. [Chain-of-thought prompting elicits reasoning in large language models](#). In *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22*, Red Hook, NY, USA. Curran Associates Inc.
- Jason Wei et al. 2023. [Chain-of-thought prompting elicits reasoning in large language models](#). *Preprint*, arXiv:2201.11903.
- Yanbin Wei et al. 2024. [Kicgpt: Large language model with knowledge in context for knowledge graph completion](#). *Preprint*, arXiv:2402.02389.
- Y. Weng, M. Zhu, G. Bao, H. Zhang, J. Wang, Y. Zhang, and L. Yang. 2025. [CycleResearcher: Improving Automated Research via Automated Review](#). arXiv preprint arXiv:2411.00816. Version posted March 8, 2025.
- Yifan Wu et al. 2024. [Chartinsights: Evaluating multimodal large language models for low-level chart question answering](#). *Preprint*, arXiv:2405.07001.
- Wenda Xu et al. 2024a. [Pride and prejudice: Llm amplifies self-bias in self-refinement](#). *Preprint*, arXiv:2402.11436.
- Yunqi Xu et al. 2024b. [Face4rag: Factual consistency evaluation for retrieval augmented generation in chinese](#). *Preprint*, arXiv:2407.01080.
- Anna Xyngkou et al. 2024. [Mindtalker: Navigating the complexities of ai-enhanced social engagement for people with early-stage dementia](#). In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. ACM.
- Lin Yang et al. 2024a. [On the evaluation of large language models in unit test generation](#). *Preprint*, arXiv:2406.18181.
- Qisen Yang et al. 2024b. [Psychogat: A novel psychological measurement paradigm through interactive fiction games with llm agents](#). *Preprint*, arXiv:2402.12326.
- Youngsik Yun and Jihie Kim. 2024. [Cic: A framework for culturally-aware image captioning](#). *Preprint*, arXiv:2402.05374.
- Chen Zhang et al. 2024a. [Teaching large language models an unseen language on the fly](#). *Preprint*, arXiv:2402.19167.
- Jintian Zhang et al. 2024b. [Exploring collaboration mechanisms for llm agents: A social psychology view](#). *Preprint*, arXiv:2310.02124.
- Yang Zhao et al. 2024. [Deciphering the impact of pre-training data on large language models through machine unlearning](#). *Preprint*, arXiv:2402.11537.

Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Lingpeng Kong, Jiajun Chen, Lei Li, and Shujian Huang. 2023. [Multilingual machine translation with large language models: Empirical results and analysis](#). *ArXiv*, abs/2304.04675.

## A Redaction Automation Attempts

We detail all prompts and specific examples for each approach in Appendix D and Appendix E, respectively. Note that all of our automation approaches start off with the filtered paper, which as a reminder is the trimmed paper but without the methodology.

### Basic Prompting

One of our initial and simplest approaches was to few-shot prompt a capable LLM like GPT-4o with a `redaction_prompt` and the filtered paper. While simple, this approach’s weakness is that it fails to capture subtle mentions of the methodology in sections like the introduction and related works and is generally inconsistent.

### Sentence-By-Sentence Redaction

To improve on our basic prompting approach, we tried a new approach which we call sentence-sentence redaction. It starts by first using NLTK’s sentence tokenizer to split the filtered paper into individual sentences and then for each individual sentence, we prompt GPT-4o with a `methodology_detection_prompt` to determine if that individual sentence contains the original methodology. While better than our first approach and majorly improving on capturing finer details, the main flaw of this approach was that focusing on a single sentence takes away important context that’s necessary for determining whether it’s a sentence that describes a methodology or not. We tried to alleviate this issue by providing each sentence with pre-context (one sentence before) and post-context (one sentence after) but the improvements were marginal and often led to false positives, thus removing crucial information necessary in the prediction phase.

### Keyword (Phrase) Redaction

Our final automated approach is something we refer to as keyword (phrase) redaction, which splits the filtered paper into *thoughts* and then checks if specific keywords or phrases are in that thought. The key hypothesis behind this approach is that all writing can be split into what we call thoughts. We

define thoughts as a collection of sentences that refers to one unifying idea. For example, the excerpt “I went to the Louvre yesterday. I had to drive my car there. I got stuck in a traffic jam but I eventually made it out. At the museum, I had lots of fun. I looked at many cool paintings and sculptures. My favorite was the Mona Lisa” can be split into “I went to the Louvre ... eventually made it out” and “At the museum, I ... was the Mona Lisa” where the first thought describes how they got to the Louvre and the second thought describes what they did at the Louvre. If we can effectively apply this idea and split research papers into these thoughts, we fix both the issue of lack of granularity and lack of context.

Additionally, since authors of papers tend to use phrases such as “we propose” or “in our paper” to describe their contributions and use specific proper nouns or acronyms for the contributions they have such as “WIPE” [Tavasoli et al. \(2024\)](#), we can simply check if a keyword or phrase is in a thought, and redact it if so. We prompt GPT-4o with a `split_prompt` to split the text into thoughts and extract proper nouns with a `proper_noun_prompt`, and we handcraft a list of commonly used phrases. We then iterate over every thought and check if any of them commonly used phrases are in the thought to determine if we redact the thought.

Unfortunately, two main problems arise. The first one is that GPT-4o’s thought extraction is wildly inconsistent, ranging from single sentences to entire paragraphs and does not align with our definition of a thought. The second issue is that GPT-4o has trouble extracting just the proper nouns that were proposed by the original authors and tends to include proper nouns mentioned in previous works like “CNN” or “Chain of Thought.” The first issue is especially difficult to solve (and the idea of thoughts may not even be well-defined) and thus we decided we would be better off manually redacting, though we encourage any future work improving this approach.

## B Statistical Significance Tests

We ran a series of statistical significance tests on the results of the LLMs in order to determine if there is a significant difference in their performance.

From visual inspection of the distributions of model performances in Figure 3, it can be seen that the distributions are not normal, and are skewed to the right. Therefore, a Kruskal-Wallis Test is

	claude	gemini	gpt4o	o3mini
claude	1.000000	0.000072	0.030039	0.921461
gemini	0.000072	1.000000	0.698303	0.019083
gpt4o	0.030039	0.698303	1.000000	1.000000
o3mini	0.921461	0.019083	1.000000	1.000000

Figure 2: P-values between each pair of model to indicate probability of no significant difference assuming null hypothesis is true. P-value of less than 0.05 indicates a significant difference between two models.

appropriate, as it does not assume anything about the normality of the distributions. Our hypotheses for this test are as follows:

Null Hypothesis ( $H_0$ ): There is no significant difference in benchmark performance between the 4 LLMs (GPT-4o, Gemini, Claude, GPT-03-Mini).

Alternative Hypothesis ( $H_A$ ): There is a significant difference in benchmark performance between at least one pair of the 4 LLMs.

After running the test using the SciPy library with a standard significance level  $\alpha$  of 0.05, our p-value was found to be 0.0001. Thus, we reject our null hypothesis that there is no significant difference between any of the models.

To determine which pairs of models have a significant difference, we used the `scikit_posthocs` library to run a post-hoc Dunn’s test with Bonferroni correction, as shown in Figure 2.

The models with a statistically significant difference are:

Claude vs Gemini: p-value = 0.0001  
 Claude vs GPT-4o: p-value = 0.0300  
 Gemini vs O3-mini: p-value = 0.0191

The models without a statistically significant difference are:

Claude vs O3-mini: p-value = 0.9215  
 Gemini vs GPT-4o: p-value = 0.6983  
 GPT-4o vs O3-mini: p-value = 1.0000

Lastly, a power analysis was performed to determine the probability of detecting a difference between the LLMs assuming one actually exists. The achieved power was estimated using an ANOVA-based approximation with Cohen’s  $f$  as the effect size measure. The power was calculated to be 0.9674, showing a high probability of detecting a difference between the LLMs.

## C Multiple Choice Question Evaluation Attempt

Our initial approach to evaluating the LLMs involved a multiple choice question format instead of a jury of LLM judges. The multiple choice options would contain perturbed summarized versions of the methodology with several major flaws in its execution, as well as a summarized version of the original methodology. The idea here was that if an LLM was able to consistently choose the correct methodology option, then it would demonstrate a strong understanding of what a rigorous research methodology entails. However, we found several issues with this approach. While an LLM that consistently chooses the correct methodology option could potentially demonstrate rigorous research ability, this process passively spoon-feeds the LLM with a well-thought out series of methodological steps. The LLM would never be pushed to actively engage in the extensive reasoning process of planning, organizing, and formulating these steps. This prevents us from accurately assessing their real-world potential as research assistants. Additionally, the generation of incorrect multiple choice options was more complicated than expected. We found that much of the incorrect options had *negative language*, which we define as subtle remarks that clearly indicate an answer choice is incorrect by indirectly stating flaws within its approach. Examples of *negative language* are depicted in the Figure below.

### MCQ Generation Prompt

Write 9 incorrect variations of the following summary of a research methodology with key details changed. Do not write anything about results or conclusions. Ensure that each variation IMPLICITLY includes a common pitfall of research such as lacking any of the following: bias mitigation, ethical soundness, use of control groups when applicable, quality data, reproducibility, correct application of statistical methods, practicality, etc. Output the 9 variations and verify they are each of a similar length.

### MCQ Generation Failure

**Correct Answer Choice:** The research paper investigates translation bias and the effectiveness of multilingual Large Language Models (LLMs) in cross-lingual claim verification using the XFACT dataset, which encompasses claims in 15 languages from five language families. **The dataset ensures a balanced representation of languages and an equal distribution across five veracity labels...**

**Incorrect Answer Choice:** The research paper explores the potential of multilingual Large Language Models (LLMs) in cross-lingual claim verification using the YFACT dataset, which includes claims in 10 languages from three language families. **The dataset lacks a balanced representation of languages and an unequal distribution across four veracity labels...**

Answer choices based off of "A Comparative Study of Translation Bias and Accuracy in Multilingual Large Language Models for Cross-Language Claim Verification" [Singhal et al. \(2024a\)](#)

## D Prompts

The `rewrite_paper_prompt` takes a parsed paper and turns it into JSON format based on a Pydantic schema.

### Rewrite Paper Prompt

You will be given a research paper in JSON format.

**Instructions:**

- You **MUST** output the original paper with **ALL** original content.
- Do **NOT** add OR remove any content (except for cleaning up any nonsensical text).
- You **MUST** preserve original section names.

**Output Format:**

Provide your response in valid JSON format with the following keys:

- `abstract`: str
- `introduction`: str
- `related_works`: List[str] (Include **ALL** of the following if they exist: Related Works, Background)
- `contributions`: List[Section] (Include **ALL** of the following if they exist: Methods, Experiments)

The `redaction_prompt` takes in a filtered paper and removes all mentions of the author's methodology.

### Redaction Prompt

You will be provided with a research paper in JSON format.

**Instructions:**

- Identify and remove portions of the following research paper that describe the methodology and experimental setup as it pertains to the study.
- Maintain the paper's structure, reformat only when necessary, and do not summarize, paraphrase, or modify text.

**Examples:**

We omit our six specific examples for brevity.

**Output Format:**

Provide your response in valid JSON format with the following keys:

- `abstract`: str
- `introduction`: str
- `related_works`: List[str] (Include **ALL** of the following if they exist: Related Works, Background)

The `methodology_detection_prompt` takes in the trimmed paper and a sentence from that paper and outputs a YES or a NO depending on if that sentence is directly related to the methodology.

### Methodology Detection Prompt

You will be given a research paper in JSON format and a sentence from the research paper.

**Instructions:**

- You are an expert research analyzer.
- Carefully analyze the provided sentence in the context of the research paper.
- Output "YES" if the sentence suggests ANY indication of the methodology/methods.
- Mark as "YES" if the sentence describes the method or how it is used, mentions results involving the method, or includes any keywords that hint at the method.
- Only output "NO" if the sentence is strictly background information (e.g., research domain or related works unrelated to the current method).

The `split_prompt` takes a filtered paper and splits it into thoughts.

### Split Prompt

You will be given a section of a research paper in JSON format.

**Instructions:**

- Split the given text into different thoughts.
- A thought is a collection of sentences that focuses on a singular idea.
- A thought usually starts with a general sentence that sets up the focus, followed by elaboration.
- A thought **MUST** contain 3-5 sentences.
- You **MUST** output the original text. Do **NOT** add or remove content.

The `proper_noun_prompt` takes a filtered paper and extracts proper nouns or acronyms that are

related to the author's methodology.

### Proper Noun Prompt

You will be given a research paper in JSON format.

#### Instructions:

- Proper nouns can take many forms: singular words, multiple words, acronyms, etc.
- If a proper noun has multiple words, also include its acronym.
- Choose AT MOST 6 proper nouns (excluding acronyms).
- If a proper noun is a commonly used keyword or too general, do NOT include it.
- If a proper noun has been mentioned in past research or works, do NOT include it.

The `outline_prompt` takes the redacted paper and generates an outline for a potential methodology.

### Outline Prompt

You will be given parts of a research paper in JSON format.

#### Instructions:

- You are an expert researcher.

#### Output Format:

Provide your response in valid JSON format with the following keys:

- `proposed_method`: str (Using the given information, first provide inspiration behind a new proposed method to address the main research problem. You should also motivate why the proposed method would work better than existing works. Then, explain how the proposed approach works, and describe all the essential steps.)
- `experimental_plan`: str (Break down EVERY single step in `proposed_method`. Every step MUST be executable. Cover ALL essential details such as the datasets, models, and metrics to be used, etc.)

The `writing_prompt` takes an outline and generates a full methodology.

### Writing Prompt

You will be given JSON formatted text.

#### Instructions:

- You will write the complete methodology section of a research paper in paragraph format.
- Use a formal and direct tone for the paper.
- Explain ALL steps logically with well-defined connections between ideas and sections.

- Elaborate heavily on the `experimental_plan`. Include extreme detail and thoroughness.
- Include specific examples with detailed explanations for further elaboration.
- Write all mathematical expressions in LaTeX.
- Be EXTREMELY verbose and thorough.
- Do NOT use excessive subsections. Instead, connect certain concepts within a section in a smooth way.

#### Output Format:

Provide your response in valid JSON format with the following key:

- `methodology`: List[Section] (The methodology section will include ALL of the following sections: Methods, Experiments.)

The `judging_prompt` takes a predicted methodology and ground truth methodology and generates an explanation an integer score from 1-10, based on a certain rubric. The rubric guidelines and trivial differences are provided after this prompt.

### Judging Prompt

You will be given the TRUE research paper contributions and the PREDICTED research paper contributions, both in JSON format.

#### Instructions:

- You are an honest and analytical judge.
- Compare how similar the PREDICTED methodology are to the TRUE methodology using the rubric provided below.
- Do NOT consider stylistic or writing choices, nor trivial details in your comparison.
- Prioritize clarity, correctness, and alignment of ideas with the research problem over the use of mathematical notation.
- Do NOT include any additional commentary, tags, or quotes outside the JSON.

#### Output Format:

Provide your response in valid JSON format with the following keys:

- `explanation`: str (A detailed rationale in paragraph format behind your judging. Avoid lists or bullet points.)
- `score`: int (An integer score from 1-10 (inclusive) based on the rubric provided below.)

#### Rubric:

Insert rubric here.

#### Trivial Differences:

Insert trivial differences here.

The rubric gives specific descriptions and examples for what each score from 1-10 looks like.

### Rubric

#### Rubric Specifications:

- Each paper will be assigned a score between 1 and 10 (inclusive), based on how similar it is to the original paper.
- When assigning a score, clearly explain which parts of the paper’s methodology are similar or different from the original.
- The following rubric will be used to assign scores:
  - **Score 1:** Methodology is vastly different.
  - **Score 2:** Methodology is significantly different but contains a few minor similarities.
  - **Score 3:** Methodology has some similarities but misses key details.
  - **Score 4:** Methodology is somewhat similar but still lacks important aspects.
  - **Score 5:** Methodology shows a relatively equal mix of similarities and differences.
  - **Score 6:** Methodology is fairly close to the original but omits or alters some details.
  - **Score 7:** Methodology closely matches the original with only minor alterations.
  - **Score 8:** Methodology is very similar with only small, noncritical differences.
  - **Score 9:** Methodology is nearly identical with extremely minor differences.
  - **Score 10:** Methodology is almost identical, with only trivial differences ignored.

#### Example Reference Paper:

We omit the example for brevity. Full examples can be found in our Github\*.

#### Examples of Score Assignments:

We omit the example for brevity. Full examples can be found in our Github\*.

\*\*Github to be provided upon acceptance

The `trivial_differences` describe certain differences between the predicted methodology and original methodology that can be ignored. Note that our original prompt contained an example for each category but we remove it for brevity.

### Trivial Differences

Below is a list of illustrative examples showing that minor or trivial differences in a predicted methodology should not be penalized when the overarching ideas and experimental approaches remain intact. These examples are drawn from various aspects of AI/ML research papers in general.

- Alternate Evaluation Metrics with the Same Ob-

jective

- Slight Variation in Data Preprocessing Techniques
- Different Hyperparameters with Similar Impact
- Minor Architectural Adjustments in Model Design
- Alternate Statistical Analysis in Ablation Studies
- Alternate Magnitudes of Values

## E Examples

As mentioned previously, we removed papers with extensive mathematical notation due to failure to parse  $\LaTeX$ . An example is shown below.

### $\LaTeX$ Parsing Failure

#### Original Text:

“Dunefesky et al. (2024) finds a computational graph through the MLP layers by training transcoders:

$$\mathbf{z}(\mathbf{h}_{\text{pre}})_i \left( \mathbf{W}_{\text{dec}}^{(A)\top} \mathbf{W}_{\text{enc}}^{(B)\top} \right)_{i,:}. \quad (1)$$

Here,  $\mathbf{W}_{\text{dec}}^{(A)\top} \mathbf{W}_{\text{enc}}^{(B)\top} \in \mathbb{R}^{|\mathcal{F}| \times |\mathcal{F}|}$  serves as a transition operator between the feature spaces of layers  $A$  and  $B$ , revealing which features in  $B$  are ancestors for the  $i$ th feature in  $A$ .

Matrices  $\mathbf{P}^{(A \rightarrow B)}$  and  $\mathbf{W}_{\text{dec}}^{(A)\top} \mathbf{W}_{\text{enc}}^{(B)\top}$  are in some sense similar.”

#### Parsed Text:

“Dunefesky et al. (2024) finds a computational graph through the MLP layers by training transcoders:

$\mathbf{z}(\mathbf{h}_{\text{pre}})_i$   
 $\mathbf{W}^{(A)\top}$   
 dec  $\mathbf{W}^{(B)\top}$   
 enc  
 $i, :$   
 (1)

Here,  $\mathbf{W}^{(A)\top}$

dec  $\mathbf{W}^{(B)\top}$

enc  $\in \mathbb{R}^{|\mathcal{F}| \times |\mathcal{F}|}$  serves as a transition operator between the feature spaces of layers  $A$  and  $B$ , revealing which features in  $B$  are ancestors for the  $i$ th feature in  $A$ .

Matrices  $\mathbf{P}(A \rightarrow B)$  and  $\mathbf{W}^{(A)\top}$

dec  $\mathbf{W}^{(B)\top}$

enc are in some sense similar.”

Excerpt taken from Laptev et al. (2025)

Additionally, we mentioned many attempts at automation briefly in Section 3 and also more detail in Appendix A. We first give an example of a correct redaction using our two-step manual redaction process then we give some specific examples

of the redaction breaking down for each automated approach. We use the paper by Cohn et al. (2024a) as our example. We color redacted text in red, false negatives in blue, and false positives in green. We omit some text with ellipsis for brevity. Also note that we start off with the filtered paper.

### Manual Redaction

"abstract": "This paper explores the use of large language models (LLMs) to score and explain short-answer assessments in K-12 science. While existing methods can score more structured math and computer science assessments, they often do not provide explanations for the scores. Our study focuses on employing GPT-4 for automated assessment in middle school Earth Science, combining ... for open-ended science assessments.",  
"introduction": "Improvements in Science, Technology, Engineering, and Mathematics (STEM) education have accelerated the shift from teaching and assessing facts to developing ... key to gaining a deep understanding of scientific phenomena (Mao et al. 2018).  
This paper develops an approach for human-in-the-loop .. of critical need in K-12 STEM instruction.",  
"related\_works": "To understand the difficulties students face when learning science, teachers need to ... Very little research has examined effective mechanisms for generating automated grading and useful formative feedback for K-12 students that are aligned with classroom learning goals. Advances in natural language processing (NLP) have produced improved automated assessment scoring approaches to support teaching ... While these methods have enjoyed varying degrees of success, a majority of these applications have targeted more structured mathematics and ... The data needed for training our models is small, imbalanced, and non-canonical in terms of syntax and semantics, all of which may impact model performance (Cohn 2020).  
This research tackles several critical issues, namely: (1) grading open-ended, short-answer questions focused on science conceptual knowledge and reasoning, (2) utilizing LLMs to ... generation process to resolve discrepancies and support the learning goals."

### Basic Prompting Failure

"abstract": "This paper explores the use of large language models (LLMs) to score and explain short-answer assessments in K-12 science. While existing methods can score more structured math and computer science assessments, they often do not provide explanations for the scores. Our study focuses on employing GPT-4 for automated assessment in middle school Earth Science, combining ... for open-ended science assessments.",  
"introduction": "Improvements in Science, Technology, Engineering, and Mathematics (STEM) education have accelerated the shift from teaching and assessing facts to developing ... key to gaining a deep understanding of scientific phenomena (Mao et al. 2018).  
This paper develops an approach for human-in-the-loop .. of critical need in K-12 STEM instruction.",  
"related\_works": "To understand the difficulties students face when learning science, teachers need to ... Very little research has examined effective mechanisms for generating automated grading and useful formative feedback for K-12 students that are aligned with classroom learning goals. Advances in natural language processing (NLP) have produced improved automated assessment scoring approaches to support teaching ... While these methods have enjoyed varying degrees of success, a majority of these applications have targeted more structured mathematics and ... The data needed for training our models is small, imbalanced, and non-canonical in terms of syntax and semantics, all of which may impact model performance (Cohn 2020).  
This research tackles several critical issues, namely: (1) grading open-ended, short-answer questions focused on science conceptual knowledge and reasoning, (2) utilizing LLMs to ... generation process to resolve discrepancies and support the learning goals."

**Explanation:** In this example, it essentially just returned the original filtered paper as the redacted paper, resulting in a lot of false negatives. This is clearly problematic as our basic prompting cannot pick up finer details.

### Sentence Sentence Redaction Failure

"abstract": "This paper explores the use of large language models (LLMs) to score and explain short-answer assessments in K-12 science. While existing methods can score more structured math and computer science assessments, they often do not provide explanations for the scores. Our study focuses on employing GPT-4 for automated assessment in middle school Earth Science, combining ... A systematic analysis of our method's pros and cons sheds light on the potential for human-in-the-loop techniques to enhance automated grading for open-ended science assessments."

"introduction": "Improvements in Science, Technology, Engineering, and Mathematics (STEM) education have accelerated the shift from teaching and assessing facts to developing ... key to gaining a deep understanding of scientific phenomena (Mao et al. 2018).

This paper develops an approach for human-in-the-loop LLM prompt engineering using in-context learning and chain-of-thought reasoning with GPT-4 to support auto-mated analysis and feedback generation for formative assessments in a middle school Earth Science curriculum. We present our approach, discuss our results, evaluate the limitations of our work, and then propose future research in this area of critical need in K-12 STEM instruction."

"related\_works": "To understand the difficulties students face when learning science, teachers need to ... which may impact model performance (Cohn 2020).

This research tackles several critical issues, namely: (1) grading open-ended, short-answer questions focused on science conceptual knowledge and reasoning, (2) utilizing LLMs to ... generation process to resolve discrepancies and support the learning goals."

**Explanation:** Our sentence-sentence approach is slightly better and has some correct redactions but fails to recognize other sentences related to the methodology, probably due to lack of context.

For our keyword approach, we separate the text in a list to represent the thoughts that it was split into, we bold key phrases like "our method", and we also include all the extracted proper nouns.

### Keyword (phrase) Redaction Failure

**Proper Nouns:** GPT-4, Chain-of-thought Reasoning, In-context Learning, Human-in-the-loop LLM Prompt Engineering, formative assessments

"abstract": ["This paper explores the use of large language models (LLMs) to score and explain short-answer assessments in K-12 science. While existing methods can score more structured math and computer science assessments, they often do not provide explanations for the scores. Our study focuses on employing GPT-4 for automated assessment in middle school Earth Science, combining ... for open-ended science assessments."],

"introduction": ["Improvements in Science, Technology, Engineering, and Mathematics (STEM) education have ... time-consuming for teachers and susceptible to errors (Rodrigues and Oliveira 2014; Haudek et al. 2011).", "Large Language Models (LLMs) provide opportunities for automating short answer scoring (Funayama et al. 2023) and ... are key to gaining a deep understanding of scientific phenomena (Mao et al. 2018).",

"This paper develops an approach for human-in-the-loop LLM prompt engineering using ... in this area of critical need in K-12 STEM instruction."],

"related\_works": ["To understand the difficulties students face when learning science, teachers ... and reasoning to better support their developing STEM ideas (Cizek and Lim 2023).",

"However, grading formative assessments, ... aligned with classroom learning goals.",

"Advances in natural language processing (NLP) have produced improved automated assessment scoring approaches to support teaching and learning (e.g., Adair et al. 2023; Wilson et al. 2021). ... different from scoring free-form short-answer responses by middle school students in science domains."

"Data impoverishment concerns ... in terms of syntax and semantics, all of which may impact model performance (Cohn 2020). This research ... aligned with specified learning objectives for both students and teachers and (3) addressing concerns related to data impoverishment. We hypothesize that our approach supports automated scoring and explanation that (1) aligns with learning objectives and standards, ... discrepancies and support the learning goals."]

**Explanation:** As you can see, there are a ton of false positives with this approach and this is due to the fact "formative assessments" was extracted as a proper



noun, which is too general. Additionally, because our thoughts were too large, this further encouraged false positives as non-methodology revealing content would be combined with the methodology.

Now, we will present some examples from the prediction pipeline that we believe is representative of the general quality of LLM generated methodologies. Specifically, for each predictor LLM, we include an example that received a score close to the average final score for that predictor LLM. We only include the `experimental_plan` of the outline for the sake of length and we also include a shortened original methodology for reference.

### GPT-4o Prediction

**Paper Citation:** Paper from [Dutta and Jannesari \(2024\)](#).

#### Original Methodology:

- Data Collection:** Compile programs from datasets (CodeNet, HPCorpus, LS-CAT) into LLVM IRs. Select 10,000 unique programs.
- Tokenization:** - Train a WordPiece tokenizer specific to the structure of IRs. - Define special tokens and use a sequence length of 64 for IR tokens.
- Graph Generation:** - Use PROGRAML to convert IRs into multi-graphs containing data-flow, control-flow, and call-flow.
- Pre-Training Tasks:** - **Masked Language Modeling:** For each IR statement: - Randomly mask 15% of tokens with specific replacements (80% '[MASK]', 10% random, 10% unchanged). - Use transformer layers to train the model to predict masked tokens. - **Graph Auto-Encoding:** - Extract sub-graphs (control-flow, data-flow, call-flow) from IRs. - Use GCN layers to encode and decode each sub-graph, optimizing graph reconstruction via adjacency matrices. - Aggregate losses for overall graph learning. - **IR-Graph Matching:** - Generate data pairs (positive from same IR, negative from different IRs). - Embed IR statements into sequences, aggregate vectors for entire files. - Train using binary cross-entropy to classify if sequence and graph match.
- Evaluation:** - Use the pre-trained model on downstream tasks to evaluate performance optimization. - Compare results against baseline methods to validate the improvement in capturing semantic and structural code features.

#### Predicted Methodology:

- Data Collection:** - Gather a large corpus of open-source HPC applications in a variety of programming languages. - Use LLVM tools to generate Intermediate Representations (IRs) for each application. - Collect dynamic profiling logs during execution on different HPC architectures to capture runtime behavior.
- Pre-processing:** - Convert LLVM IR into multi-graph representations using a tool like PROGRAML. - Extract static features from IRs and dynamic features from execution logs. - Normalize and clean the dataset to ensure consistency across

different programs and architectures.

- Multi-modal Representation Learning:** - Train separate encoders for graph-based LLVM IR representations and dynamic profiles. Use Graph Neural Networks (GNNs) for IR graphs and recurrent neural networks (RNNs) for dynamic profiles. - Develop a fusion layer that combines the learned embeddings from both encoders to form a unified representation.
- Pre-training HPC-BERT:** - Pre-train the combined model using a self-supervised learning approach where the model learns to predict missing parts of code or performance metrics. - Employ masked node prediction tasks in the IR graph and masked performance prediction tasks in the dynamic profile. - Use a large-scale dataset to ensure robust learning of code and performance dependencies.
- Fine-tuning for Downstream Tasks:** - Fine-tune HPC-BERT on specific HPC optimization tasks such as auto-tuning compiler flags or predicting parallelization strategies. - Use datasets specific to each downstream task, ensuring alignment with the pre-training phase. - Evaluate performance improvements against baseline ML methods and human-tuned benchmarks using metrics like execution time reduction and energy efficiency.
- Evaluation and Validation:** - Conduct comprehensive experiments on HPC clusters to assess the model's adaptability and performance gains. - Validate the model's predictions against human expert optimizations and iterate on the model based on feedback and results.

**Score:** 3.33

### o3-mini Prediction

**Paper Citation:** Paper from [Rasool et al. \(2024\)](#).

#### Original Methodology:

- Query and Answer Synthesis:** - Obtain a corpus of documents that are semantically represented (i.e., embedding vectors). - Use a large language model (LLM) to extract factual statements from each document. - For each extracted fact, employ the LLM to generate questions that use these facts as answers, ensuring the questions remain within the document's scope. - If dealing with long or multiple documents, use a retriever to extract relevant facts for question generation.
- Query Evaluation:** - For each generated query, retrieve the top N documents and check if the original document is among them. - Retain queries that successfully retrieve relevant documents and discard those that do not, ensuring domain-specific relevance.
- Query Variation Generation:** - Prompt the LLM to create diverse variations of each verified query, focusing on differences in phrasing, structure, or keywords while retaining semantic similarity. - Provide both the question and answers in the prompt to maintain domain relevance and mitigate unrelated or hallucinated queries.
- Evaluation Strategy:** - Use the Qasper dataset, which includes question-answers from NLP research papers, filtering out unanswerable questions. - Utilize the Assignment dataset, a curated set of 52 distinct questions derived from typical university-level queries. - Perform qualitative evaluations involving human assessment of generated question-answer

pairs for realism and relevance. - Conduct a case study using the Assignment dataset with the GPT-Cache semantic caching library, evaluating performance with DistilBERT, SBERT, and ONNX embeddings. - Measure metrics such as correct/incorrect cache hits and misses to assess the semantic cache's effectiveness with the generated dataset.

#### **Predicted Methodology:**

**\*\*Step 1: Dataset Preparation\*\*** - Gather a small set of domain-specific documents and existing query-response pairs to fine-tune the language model. This will also serve as the basis for training the GAN discriminator.

**\*\*Step 2: GAN Model Training\*\*** - Initialize a GAN with a generator capable of creating varied queries and a discriminator pre-trained on a general language model and further fine-tuned using the domain-specific dataset. - Train the GAN to produce semantically diverse and domain-relevant queries, iteratively improving the generator's performance based on discriminator feedback.

**\*\*Step 3: Fine-Tuning Language Model\*\*** - Fine-tune a large language model (e.g., BERT or GPT) with the domain-specific dataset. This model will serve as the discriminator within the GAN framework to ensure relevance and domain consistency of generated queries.

**\*\*Step 4: Adversarial Query Generation\*\*** - Use the trained GAN to generate a wide range of queries, including slight perturbations that act as adversarial examples. These queries should vary in semantic proximity to existing ones to test cache robustness.

**\*\*Step 5: Implementation of Semantic Cache Testing\*\*** - Deploy the newly created queries against the semantic cache system. - Record data on cache hits and misses, focusing on how adversarial and domain-specific queries are handled.

**\*\*Step 6: Feedback Loop Integration\*\*** - Analyze incorrect cache hits and misses, feeding this data back into the GAN training process. - Modify the generator and discriminator based on these insights to refine their outputs and improve future query generation.

**\*\*Step 7: Evaluation Metrics and Iteration\*\*** - Set evaluation metrics, such as precision, recall, and F1-score, to measure the performance of cache hits/misses. - Continuously iterate on the GAN model based on evaluation results to enhance query generation quality.

**\*\*Step 8: Scalability Testing\*\*** - Test the scalability of the method by deploying it in various domains, adjusting the fine-tuning process as necessary for each domain's specific characteristics.

**Score:** 3.33

### **Claude 3.5 Sonnet Prediction**

**Paper Citation:** Paper from Salminen et al. (2024).

#### **Original Methodology:**

- \*\*Select the Research Context\*\*:** Focus on the domain of addiction, covering five types: alcohol, opioids, social media, online shopping, and gambling.
- \*\*Use GPT-4 for Persona Generation\*\*:** - Set up the experiment with three types of prompts for each addiction: male specified, female specified, and gender-neutral. - Generate 450 personas by repeating the generation 30 times for each of the 15 prompt

combinations using GPT-4.

#### **3. \*\*Implement a Two-Stage Prompting Strategy\*\*:**

- Stage 1: Generate skeletal personas for each addiction-prompt combination. - Stage 2: Ask GPT-4 to expand the skeletal descriptions into full personas.

**4. \*\*Ensure Replicability\*\*:** - Provide a structured template for persona information output. - Share the generated Personas-addicted dataset for further research.

**5. \*\*Internal Evaluation\*\*:** - Conduct an internal evaluation with four experienced UX researchers to assess persona quality using both quantitative and qualitative metrics. - Compute inter-coder reliability using Krippendorff's Alpha (U) metric.

**6. \*\*External Evaluation\*\*:** - Recruit five subject-matter experts (SMEs) in public health via Upwork to evaluate a random stratified sample of 30 personas. - Use a standardized evaluation framework, covering criteria like age, gender, occupation, and personality.

**7. \*\*Analyze Results\*\*:** - Conduct statistical tests (e.g., Chi-squared, Mann-Whitney U) to identify biases in gender, age, and country distributions. - Use regression modeling to explore relationships between pain points and demographic variables.

**8. \*\*Provide Feedback and Adjustments\*\*:** - Collect qualitative feedback from SMEs on persona realism, relatability, and usability. - Use feedback to refine prompt strategies and mitigate biases.

**9. \*\*Share Findings\*\*:** - Publish the methodology and results, highlighting the potential of LLMs for persona generation and addressing biases and quality feedback.

#### **Predicted Methodology:**

**1. \*\*Dataset Preparation\*\*:** Collect a comprehensive dataset comprising various demographic statistics from trustworthy sources, such as national census data, to serve as benchmarks for diversity in personas.

**2. \*\*LLM Persona Generation\*\*:** Use a pre-trained large language model like OpenAI's GPT-4 to generate initial personas based on prompts describing different user groups.

**3. \*\*Bias Detection\*\*:** Implement a bias detection algorithm using natural language processing techniques and fairness-aware machine learning models. This step involves analyzing the generated personas to identify potential demographic imbalances or biased representations.

**4. \*\*Diversity Benchmarking\*\*:** Compare the demographic attributes of the generated personas against real-world statistics collected in step 1. This comparison helps in identifying under-represented or over-represented groups.

**5. \*\*Fairness-Aware Filtering\*\*:** Apply a filtering mechanism to adjust the features of the personas to better match the diversity benchmarks. This can involve tweaking attributes such as ethnicity, age, gender, and other relevant demographic factors to ensure a more balanced representation.

**6. \*\*Iterative Refinement\*\*:** Re-run the generation and filtering processes iteratively, refining the prompts and fairness criteria based on feedback from domain experts until the personas meet acceptable levels of diversity and bias reduction.

**7. \*\*Evaluation\*\*:** Conduct a comprehensive evaluation with subject-matter experts to assess the improved personas for engagement, interestingness, realism, diversity, and bias. Use metrics like the diversity index and sentiment analysis to quantify im-

provements.

8. **User Perception Study**: Perform user studies to gather feedback on the perceived realism and inclusiveness of the personas. This feedback loop helps in further refining the generation process.

9. **Documentation and Reporting**: Document the entire process, detailing the methods, findings, and improvements made to the LLM-generated personas. Provide recommendations for future research and potential applications of PersonaFair in various domains.

**Score:** 2.67

## Gemini 1.5 Pro Prediction

**Paper Citation:** Paper from Toki et al. (2024).

### Original Methodology:

1. **Data Collection:** Recruit participants, primarily young children, through health and educational establishments. Obtain informed consent from parents and record the children's developmental and communication skills. Collect data using the SmartSpeech system's serious game activities designed to gather linguistic and biometric measurements.

2. **Data Preprocessing:** - Handle missing data using a model-based imputer in Orange data mining software. - Select 13 relevant features. - Standardize variables to mean = 0, SD = 1. - Filter out outliers to maintain data integrity.

3. **Cluster Analysis:** - Conduct hierarchical clustering using Spearman distance metric and Ward linkage to identify clusters of similar attributes. - Analyze the resulting dendrogram to understand the variable groupings.

4. **Reliability Analysis:** - Calculate Cronbach's alpha for each cluster to assess internal consistency and reliability.

5. **Factor Analysis:** - Perform exploratory factor analysis using minimum residual extraction and varimax rotation to identify latent factors. - Confirm factor structure using Bartlett's test of sphericity and KMO measure.

6. **Predictive Modeling:** - Develop a logistic regression model using the latent factors as input features. - Use Orange Data Mining for model development and validation. - Evaluate model performance with metrics like accuracy, precision, recall, and F1 score. Apply stratified 10-fold cross-validation. - Perform external validation with a cohort of 184 participants, ensuring demographic representation.

7. **Evaluation:** - Analyze model performance using confusion matrix and ROC curves to assess discriminative power. - Interpret results using a nomogram for clinical insights. - Compare with alternative models (e.g., random forest, SVM) to ensure best model selection.

8. **Documentation and Reporting:** - Produce detailed reports and visualizations of model performance and statistical analyses. - Prepare for publication and dissemination of findings in clinical and academic settings.

### Predicted Methodology:

1. **Data Collection Infrastructure**: - Set up a controlled environment with audio recording devices, eye-tracking systems, and motion sensors. - Use a standardized protocol to ensure consistency across all

data collection sessions. - Obtain ethical approval and informed consent from participants' guardians. - Recruit a diverse sample of children, ensuring representation across different neurodevelopmental disorders for the training dataset.

2. **Data Preprocessing**: - For speech data, employ noise reduction filters and segment audio clips into meaningful units (e.g., words, phrases). - Apply MFCC to extract fundamental speech features. - Analyze eye-tracking data to identify and log fixations, saccades, and blink rates. - Process motor data to quantify movement patterns using algorithms like dynamic time warping for temporal alignment.

3. **Feature Extraction and Selection**: - Use Python libraries such as librosa for audio features, pandas and numpy for data manipulation, and OpenCV for processing eye-tracking data. - Implement feature selection techniques like recursive feature elimination to identify the most informative features.

4. **Model Development - Multi-Modal Fusion**: - Design a neural network architecture with branches for each modality, converging into a fusion layer. - Use frameworks like TensorFlow or PyTorch to implement deep learning models. - Initialize with pre-trained models where applicable to enhance learning efficiency.

5. **Training and Validation**: - Split the dataset into training, validation, and test sets (e.g., 70-15-15 split). - Use stratified sampling to ensure balanced representation of classes. - Apply data augmentation techniques to increase variability and robustness. - Employ metrics such as accuracy, F1 score, and confusion matrix for evaluation.

6. **Real-Time System Development**: - Develop a user-friendly interface using frameworks like Django or Flask for clinicians to interact with the system. - Implement visualization tools using libraries like matplotlib and seaborn to display multi-modal data insights.

7. **Iterative Feedback Loop**: - Set up a mechanism for clinicians to provide feedback on the system's outputs. - Use this feedback to iteratively refine model parameters and improve prediction accuracy.

**Score:** 3.75

In addition to examples of predictions, we will also include three examples of the judging outputs for each judge LLM. In order to demonstrate the full extent of our judging criteria, we include a judging response that's rated high (7-10), medium (4-6), and low (1-3). To be concise, we do not include all the judging responses during the voting, only the voted upon response. Also, we arbitrarily select examples from the judging results of any of the predictor LLMs.

## GPT-4o Judging (Low Rating)

**Paper Citation:** Paper from Saxena et al. (2024).

### Explanation:

The predicted methodology introduces a novel framework called ConsistentReason, which significantly

deviates from the original methodology. While both papers focus on evaluating consistency and reasoning capabilities of LLMs, the approaches are fundamentally different. The original methodology uses the Boolq dataset and focuses on evaluating LLMs in a zero-shot setting without any additional frameworks or model modifications. It assesses consistency by repeating queries and measuring response stability, and reasoning by comparing generated explanations with ground truth using metrics like BERT Score, BLEU Score, and F-1 Score. In contrast, the predicted methodology proposes a complex multi-component framework involving knowledge graph integration, iterative self-verification, and reasoning chain validation to improve factual accuracy and consistency. It uses different datasets and evaluation metrics, such as Hallucination Rate and Consistency Score, and employs a sophisticated experimental setup including human evaluations and cross-validation. These substantial differences in approach, datasets, evaluation metrics, and experimental setup justify a score of 1, as the predicted methodology is vastly different from the original.

**Score:** 1

### GPT-4o Judging (Medium Rating)

**Paper Citation:** Paper from [Liu et al. \(2024b\)](#).

#### **Explanation:**

The predicted methodology for the Dynamic Scaffolding Framework (DSF) shares some foundational similarities with the true contribution on Multi-Modal Tutoring Systems but diverges in several key areas, warranting a score of 5. Both methodologies focus on language learning through the integration of multimodal systems and the use of scaffolding strategies, which is a critical overlap. The true methodology uses GPT-4V for its image description tutoring, emphasizing pedagogical strategies like role and task definition, while the predicted methodology uses GPT-4 fine-tuned on teaching dialogues, incorporating explicit scaffolding strategies and implicit contextual cues. This difference in the LLM's application suggests a divergence in the approach to leveraging language models. Furthermore, the predicted methodology introduces a dynamic fading mechanism for scaffolding reduction, a concept not present in the true methodology, which focuses more on the pedagogical theories and scaffolding strategies like dialogic teaching, inquiry-based learning, and ZPD. The evaluation in the true paper involves scaffolding strategies and their impact on students' language learning, whereas the predicted methodology employs a rigorous experimental design comparing DSF and traditional ITS, focusing on engagement metrics, learning efficiency, and user satisfaction. These differences in evaluation approaches and the absence of specific pedagogical strategy testing in the predicted methodology indicate a balanced mix of similarities and differences, justifying a score of 5.

**Score:** 5

### GPT-4o Judging (High Rating)

**Paper Citation:** Paper from [Long et al. \(2024\)](#).

#### **Explanation:**

The predicted methodology closely aligns with the true methodology in many key aspects, though there are some notable differences. Both approaches involve the concept of generating multiple expert perspectives to address an input query, and both utilize a structured process of aggregating these perspectives to create a final response. The predicted methodology mirrors the true paper's two-step process: generating expert identities and responses, and then aggregating these responses. However, there are some deviations in execution and detail. The predicted approach describes the use of a similarity function and a ranking or voting algorithm to aggregate responses, which differs from the true methodology's detailed seven-step aggregation process that involves specific subtasks like identifying agreed viewpoints, resolving conflicts, and generating isolated viewpoints. Additionally, the predicted methodology includes an analytical module for generating expert profiles, which is not explicitly mentioned in the true paper. Despite these differences, the core principles and goals of leveraging diverse expert opinions to improve response quality and mitigate bias are consistently represented. These similarities, combined with some differences in the technical execution of aggregation and expert generation, warrant a score of 7.

**Score:** 7

### o3-mini Judging (Low Rating)

**Paper Citation:** Paper from [Lee and Lim \(2024\)](#).

#### **Explanation:**

The true research paper contributions are focused on diagnosing and analyzing the limitations of language models with respect to their lack of sensory experience – specifically through tasks such as H-TEST and Letter Geometry – in order to demonstrate that sensory aspects like visual and auditory information are blind spots in standard LLMs. This work is essentially an empirical and analytical study that uses task-based experiments and ablation studies to reveal these limitations, along with observations on few-shot prompting and chain-of-thought effects. In contrast, the predicted methodology outlines the design and implementation of a new multi-modal model architecture (SELM) that explicitly integrates multiple sensory modalities (vision, audio, and haptic) into a BERT-large based framework using dedicated sensory streams, fusion through cross-attention, and a multi-component loss function. Moreover, the experiments in the predicted methodology are focused on training a model with multi-modal data along with standard NLP benchmarks and novel sensory consistency metrics. This approach is aimed at building a new model architecture rather than evaluating the current shortcomings of sensory-deprived LLMs. As a result, the overarching objectives, experimental designs, and intended contributions are vastly different, with little overlap in the core methodological approach.

**Score:** 1

### o3-mini Judging (Medium Rating)

**Paper Citation:** Paper from [Cambrin et al. \(2024\)](#).

**Explanation:**

The predicted methodology shares several high-level ideas with the original paper. Both works propose a compound loss framework that starts with cross-entropy and integrates additional losses (Dice and Focal) to improve token prediction and address class imbalance. They both focus on training large language models using parameter-efficient methods such as LoRA and target mathematical reasoning tasks along with general language modeling. However, the predicted methodology diverges in several important respects. It omits key loss functions such as Lovász loss and self-adjusting Dice loss that are central in the original paper, and it introduces a novel token-level attention mechanism to weight the loss components—a detail that is not present in the original description. In addition, while both works perform ablation studies, the predicted experiments include varied curriculum learning strategies, different datasets (e.g., WikiText-103, C4, and Stack Exchange for general language modeling versus the original’s mixture of QA and math-specific datasets) and evaluation metrics (using ROUGE-L, BLEU, and perplexity versus Exact Match, IoU, and ROSCOE metrics). These differences in the choice of losses, additional architectural innovations, and experimental design present a balance of similarities and departures. Overall, the predicted methodology exhibits a relatively equal mix of similar elements and significant differences from the original, warranting a score that reflects this balanced alignment.

**Score:** 5

### Claude 3.5 Sonnet Judging (Low Rating)

**Paper Citation:** Paper from [Sharma et al. \(2024b\)](#).

**Explanation:**

The predicted methodology differs significantly from the true methodology in both approach and objectives. The true paper focuses on analyzing and comparing human-generated summaries from different modalities (speech vs. text) and annotator types (expert vs. non-expert), while the predicted methodology describes developing a multi-modal summarization system with contrastive learning and human-in-the-loop feedback. The true paper’s methodology centers around data collection, evaluation metrics, and research questions examining human summarization behavior, whereas the predicted paper proposes an AI system for automated summarization. While both papers deal with speech and text modalities in summarization, their fundamental approaches and goals are vastly different. The true paper is an analytical study of human summarization behavior with carefully designed data collection procedures and evaluation metrics to compare summaries across modalities and annotator expertise levels. In contrast, the predicted methodology focuses on building a deep learning system with specific architectural components like multi-modal transformers and contrastive learning. The evaluation approaches also differ substantially - the true paper uses metrics to compare human-generated summaries, while the predicted paper suggests evaluating an AI system’s output. These fundamental differences in approach, objectives, and evaluation make the methodologies largely incompatible.

**Score:** 1

### o3-mini Judging (High Rating)

**Paper Citation:** Paper from [Li et al. \(2024c\)](#).

**Explanation:**

The predicted paper’s methodology is nearly identical in structure and intent to the original work. Both methods follow a staged approach that begins with pre-training a base language model on a large, unlabeled corpus, followed by the integration of adapter modules and the use of a meta-learning (MAML) framework on source data, and concluding with a domain-adaptive fine-tuning phase on limited target data. The predicted description mirrors the original by emphasizing the utility of adapter modules to preserve general language understanding while rapidly adapting to low-resource settings, and both include comprehensive evaluation strategies using quantitative metrics. The differences are minor – for example, the specific model (BART) is named in the original but omitted in the predicted version, and there is a slight variation in the evaluation metrics (using METEOR instead of iBLEU). These details are considered trivial since they do not affect the overarching ideas and experimental framework. Thus, the core methodology, multi-stage approach, and experimental validation are extremely well-aligned between the two.

**Score:** 9

### Claude 3.5 Sonnet Judging (Medium Rating)

**Paper Citation:** Paper from [Li et al. \(2024b\)](#).

**Explanation:**

The predicted methodology shares some key elements with the original paper but has notable differences in both approach and evaluation methods. The original paper focuses heavily on using existing large language models (LLMs) with carefully crafted prompts for text style transfer across multiple datasets (YELP, CAPTIONS, GENDER, POLITICAL), employing both automated metrics and a novel duel-ranking human evaluation method. While the predicted methodology also proposes using LLMs and incorporating human evaluation, there are significant differences in the specific approaches. The original paper develops a specific prompt template with four core components and uses established automated metrics (BLEU, BERTscore, FastText classifier, perplexity) alongside their novel duel-ranking human evaluation. In contrast, the predicted methodology proposes a more general ‘Adaptive Style Transfer with Uniform Evaluation’ (ASTUE) framework with adaptive contextual learning and integrated feedback loops, which weren’t part of the original approach. The predicted methodology’s evaluation framework, while also incorporating human evaluation, lacks the specific duel-ranking mechanism that was central to the original paper’s contribution. However, both methodologies share the core idea of combining automated and human evaluation metrics, and both emphasize the im-

portance of standardized evaluation protocols. The experimental approach in the predicted methodology is more generalized and lacks the specific datasets and metrics used in the original paper, though it maintains similar high-level goals of assessing content preservation, style transfer strength, and fluency.

**Score:** 5

### Claude 3.5 Sonnet Judging (High Rating)

**Paper Citation:** Paper from [Li et al. \(2024c\)](#).

**Explanation:**

The predicted methodology shows strong alignment with the original paper's core approach while maintaining most key components, but introduces some variations in the presentation and organization of ideas. Both methodologies focus on a three-stage learning paradigm involving pre-training, meta-learning with adapters, and fine-tuning. The key similarities include: using BART as the backbone model, incorporating adapter layers into transformer architectures, employing MAML for meta-learning, and following the same overall workflow from pre-training through adaptation to target domain fine-tuning. The predicted methodology correctly emphasizes the importance of preserving prior knowledge while adapting to new domains and maintains the focus on low-resource scenarios. However, there are some differences in how certain aspects are presented and elaborated. The predicted methodology spends more time discussing data collection and evaluation strategies, while the original is more concise and mathematically formal in its presentation. The predicted version also expands more on the iterative improvement process, which wasn't explicitly detailed in the original. While these differences represent expansions rather than contradictions, they do represent slight deviations from the original presentation. Additionally, while the predicted version captures the essence of the adapter architecture, it doesn't provide the same level of mathematical detail in describing the adapter function. Despite these minor variations, the core methodological approach remains fundamentally aligned with the original paper.

**Score:** 8

### Grok 2 Judging (Low Rating)

**Paper Citation:** Paper from [Lee and Lim \(2024\)](#).

**Explanation:**

The predicted contributions from the SELM paper focus on developing a novel architecture, the Sensory-Enriched Language Model (SELM), which integrates multi-modal sensory information into language model pre-training. This approach fundamentally differs from the true contributions of the H-TEST paper, which investigates the limitations of language models in understanding visual and auditory aspects of language without sensory inputs. The SELM's methodology involves constructing a model with specialized sensory processing components (e.g., ResNet50 for vision, Wav2Vec for audio, a custom CNN for haptic), which is entirely absent from the H-TEST's approach that solely uses language-only models. The SELM paper also discusses a detailed

multi-component loss function and a sensory fusion layer, whereas the H-TEST paper does not engage in model development but rather conducts classification tasks to test the sensory understanding of existing models. Furthermore, the experimental setup in the SELM paper involves extensive dataset preparation, multi-modal training, and evaluation using standard and novel metrics, which contrasts sharply with the H-TEST's use of few-shot prompting and a focus on performance relative to a random baseline. While both papers touch on the theme of sensory understanding in language models, the SELM paper's approach to achieving this through model architecture and training is fundamentally different from the H-TEST's approach of assessing existing models' limitations. The SELM's inclusion of human evaluation and detailed ablation studies also diverges from the H-TEST's methodology, which does not explore model architecture changes or human evaluation but rather focuses on empirical testing of language model capabilities.

**Score:** 1

### Grok 2 Judging (Medium Rating)

**Paper Citation:** Paper from [Weerasinghe et al. \(2024\)](#).

**Explanation:**

The predicted contributions outline a multimodal edge-based cognitive assistant system for EMS protocol guidance, which shares several conceptual similarities with the true contributions, but also deviates significantly in its approach and execution. The true contributions detail a real-time cognitive assistant pipeline for EMS responders, focusing on speech recognition, protocol selection, and intervention recognition, all optimized for edge deployment. In contrast, the predicted contributions emphasize a comprehensive data collection process, a knowledge graph for domain knowledge representation, and the use of specific models like MobileViT v2 and Whisper-tiny, which are not mentioned in the true contributions. The predicted methodology also includes a multimodal fusion mechanism and a knowledge graph attention network, which are not part of the true contributions. While both methodologies aim to assist EMS responders, the predicted contributions focus more on technical aspects like model optimization for edge devices and detailed performance metrics, whereas the true contributions emphasize practical implementation and real-time performance. The predicted contributions also include a user study, which is absent from the true contributions, indicating a different focus on validating the system's utility. Overall, while the overarching goal of assisting EMS responders is shared, the predicted contributions diverge significantly in their methodological details and evaluation approach, warranting a score of 4 due to the lack of alignment in key aspects of the methodology.

**Score:** 4

## Grok 2 Judging (High Rating)

**Paper Citation:** Paper from Li et al. (2024c).

**Explanation:**

The predicted contributions closely align with the true contributions in terms of the overall methodology for enhancing paraphrase generation in low-resource domains. Both methodologies utilize a three-stage approach: pre-training on a large unlabeled corpus, meta-training with a meta-learning framework, and domain-adaptive fine-tuning on a target corpus. The backbone model in both is a pre-trained language model, specifically BART in the true contributions, which the predicted contributions refer to as a base language model. The use of adapter modules in the predicted contributions matches the adapter model described in the true contributions, where adapters are integrated into the transformer layers of the backbone model to facilitate domain-specific adaptation without altering the pre-trained model's core parameters. The meta-learning phase in both methodologies employs the Model-Agnostic Meta-Learning (MAML) framework, aimed at enabling rapid adaptation to new tasks with minimal data, which is crucial for low-resource scenarios. The domain-adaptive fine-tuning phase in both involves fine-tuning on a small target dataset, with the predicted contributions explicitly mentioning the use of labeled target domain data, aligning with the true contributions' use of a target corpus. The evaluation metrics listed in the predicted contributions (BLEU, ROUGE, METEOR) are similar to those used in the true contributions (BLEU, iBLEU, ROUGE), indicating a close match in performance assessment. The predicted contributions' mention of an iterative approach based on evaluation feedback is also reflected in the true contributions' emphasis on continuous improvement and adaptation to the target task. The main difference lies in the level of detail and specificity; the true contributions provide more detailed mathematical formulations and specific dataset names, whereas the predicted contributions are more general in their descriptions. However, the core ideas and experimental approaches are substantially similar, warranting a high score.

**Score:** 8

## F Supporting Figures

We present our judging pipeline in Figure 3. We present our manual redaction pipeline in Figure 4. We visualize our paper curation process in Figure 5. We include box plots to visualize the distribution of individual judge LLM scores for each predictor LLM in Figure 6. We visualize the length bias scatter plots in Figure 7.

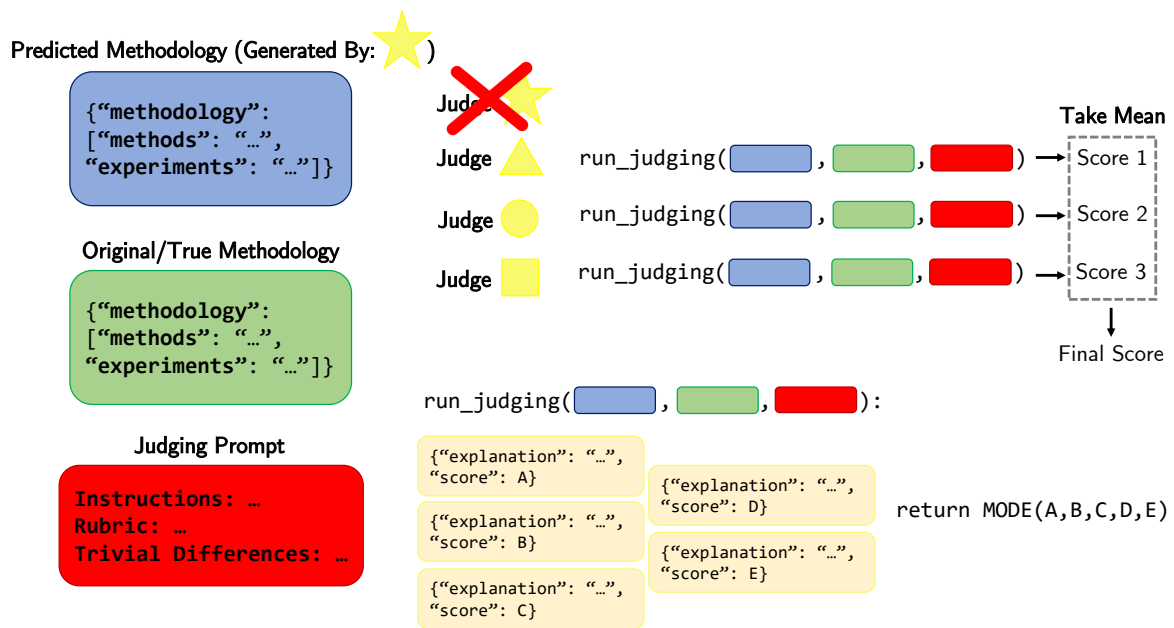


Figure 3: Visual representation of our judging pipeline with a jury of LLMs, majority voting, and self-omission.

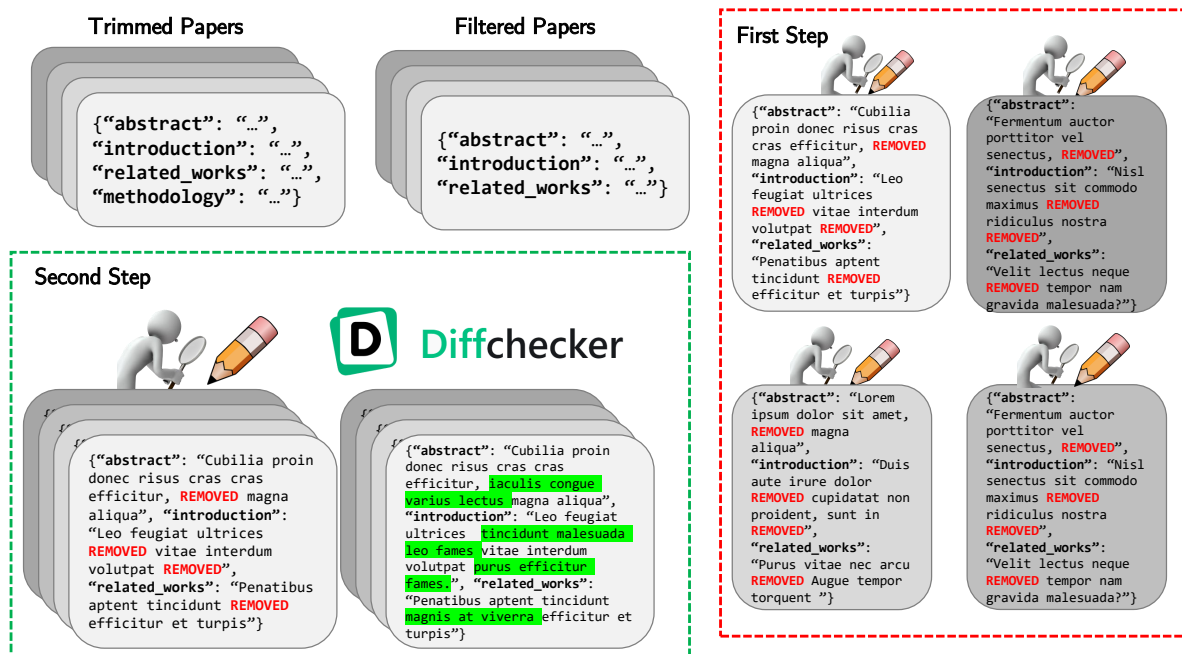


Figure 4: Visual representation of the manual redaction pipeline in two steps. Identifying information (author names, emails, affiliations) are removed along with table data, references, and appendices.



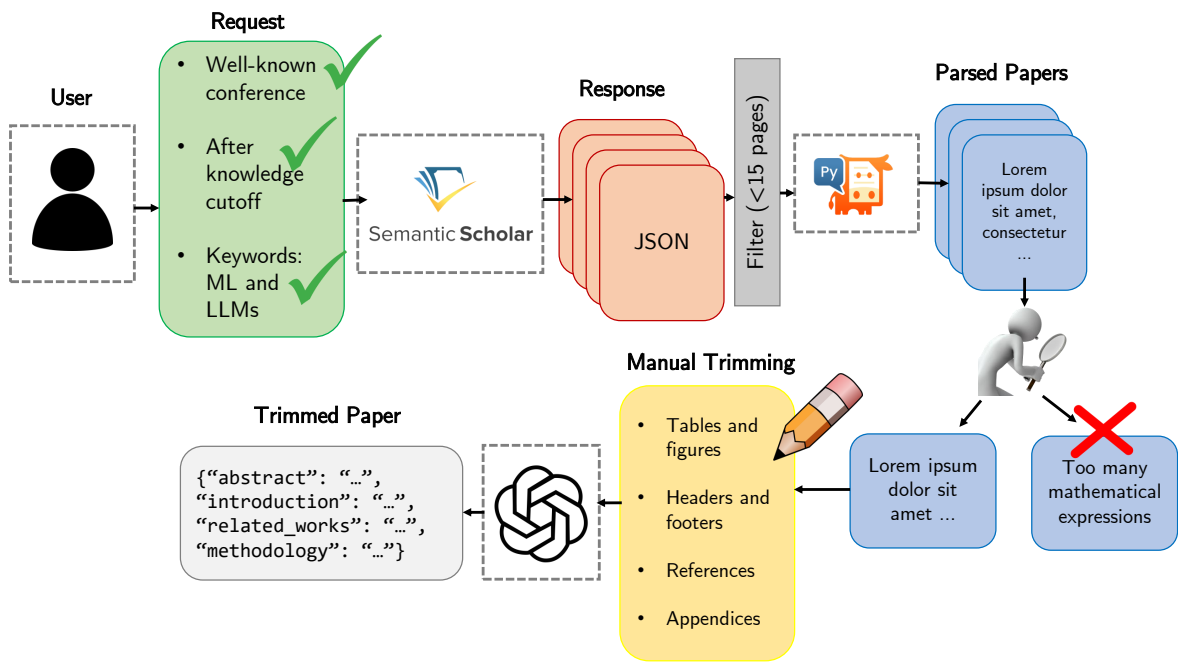


Figure 5: A visual overview of our paper curation process.

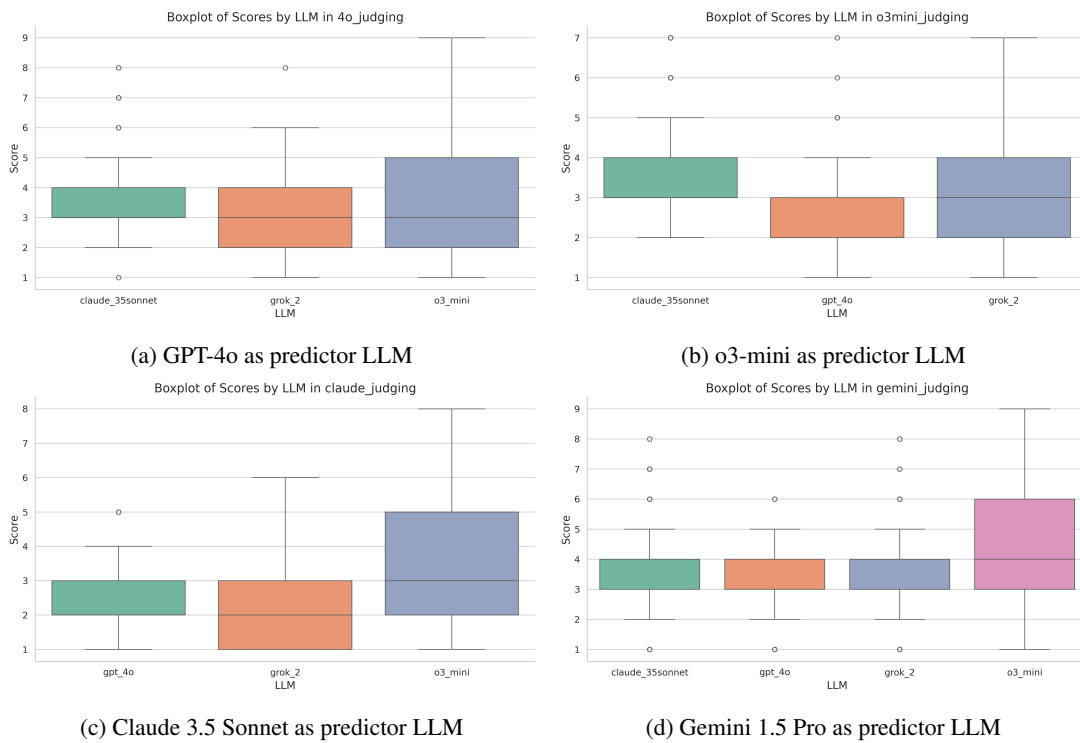


Figure 6: Box plots describing the distribution of each judge LLM’s scores that were involved for each predictor LLM. Most judge LLMs tend to give consistently low scores.

Length vs. Score (All Models)

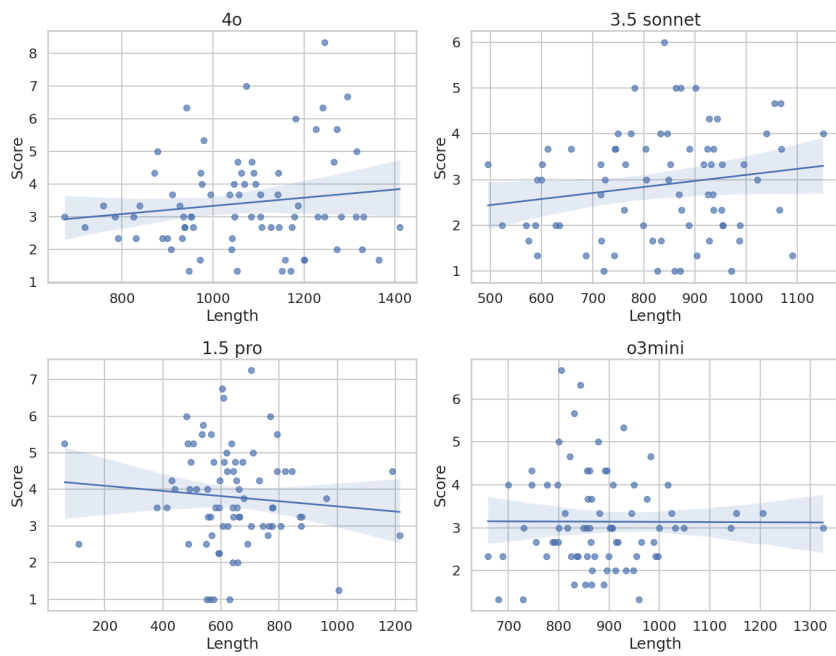


Figure 7: Scatterplot of prediction length against its score. The nearly horizontal slopes suggest there is no relationship between the length of a prediction and its assigned score.