BAGELS: Benchmarking the Automated Generation and Extraction of Limitations from Scholarly Text

Ibrahim Al Azher[†], Miftahul Jannat Mokarrama[†], Zhishuai Guo[†], Sagnik Ray Choudhury[‡], Hamed Alhoori[†]

[†]Northern Illinois University, Dekalb, IL, USA

[‡]University of North Texas, Denton, TX, USA

{iazher, mmokarrama1, zguo, alhoori}@niu.edu, sagnik.raychoudhury@unt.edu

Abstract

In scientific research, "limitations" refer to the shortcomings, constraints, or weaknesses within a study. A transparent reporting of such limitations can enhance the quality and reproducibility of research and improve public trust in science. However, authors often a) underreport them in the papers' text and b) use hedging strategies to satisfy editorial requirements at the cost of readers' clarity and confidence. This underreporting behavior, along with an explosion in the number of publications, has created a pressing need to automatically extract/generate such limitations from scholarly papers. In that direction, we report a complete architecture for computational analysis of research limitations. Specifically, we a) create a dataset of limitations in ACL, NeurIPS, and PeerJ papers by extracting them from papers' text and integrating them with external reviews; b) propose methods to automatically generate them using a novel Retrieval Augmented Generation (RAG) technique; c) create a fine-grained evaluation framework for generated limitations and provide a meta-evaluation for the proposed evaluation techniques. Code and datasets are here: github | huggingface

1 Introduction

In the context of a scientific article, "limitations" refer to the inherent shortcomings, constraints, or weaknesses within a study that may influence the results or limit the generalizability of its findings (Ross and Bibler Zaidi, 2019). These limitations can arise from various elements of the research process, including the methodology, theoretical framework, data collection, experimentation, and analysis (Ioannidis, 2007). Limitations explicitly acknowledged by the authors often include internal validity issues, measurement errors, potential confounding factors, and the failure to measure important variables (Puhan et al., 2009).

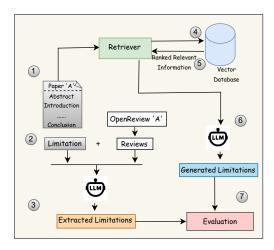


Figure 1: System Architecture. We build a dataset by extracting limitations from a paper's text (when mentioned) and reviews, build a generator for papers that do not explicitly mention limitations, and develop a novel framework for reference-based evaluation of generated limitations.

Openly discussing limitations is crucial for several reasons. It ensures the credibility and ethical standards of a scientific inquiry (Bunniss and Kelly, 2010; Chasan-Taber, 2014; Annesley, 2010). By clearly defining the boundaries of a study, researchers provide a realistic context for their findings, supporting accurate interpretation, transferability, and reproducibility (Ioannidis, 2007; Eva and Lingard, 2008). Acknowledging limitations fosters scientific integrity by demonstrating a commitment to honesty (Žydžiūnaitė, 2018). It also opens pathways for refining research methods, opening areas for future work (Azher et al., 2025), and leading to innovations that help future researchers avoid repeating the same shortcomings (Escande et al., 2016).

Despite the benefit of honest admission of such limitations, researchers sometimes remain reluctant to include them in the text or articulate them in detail (Ioannidis, 2007; Ter Riet et al., 2013).

Concerns about the potential impact on publication chances and, subsequently, career progression (Montori et al., 2004) can encourage this biased behavior. When forced to explicitly acknowledge limitations in the text¹, authors sometimes resort to using generic or irrelevant terms, which do not meaningfully inform the reader about the study's real constraints (Ross and Bibler Zaidi, 2019). Moreover, limitations can sometimes serve as a form of "hedging", where authors present their findings cautiously to avoid making definitive statements (Hyland, 1998). This approach, while safer for the authors, does not contribute to the clarity or usefulness of the research.

Failure to disclose limitations undermines the scientific process and misleads readers, reviewers, and policymakers, potentially preventing them from recognizing the constrained nature of findings and possible biases (Greener, 2018). NLP research in automatic extraction, generation, and evaluation of limitations from scholarly papers is limited due to the lack of standardized datasets, novel methods, and evaluation framework. This study takes a step toward bridging this gap.

The number of scientific publications has exploded recently (Larsen and Von Ins, 2010). That, combined with the problems of self-reporting, has necessitated a need for computational methods to study the limitations of scholarly work. In that direction, our contributions are as follows (see Figure 1):

- We create a dataset for limitations reported in research articles. For each paper, we extract the limitations from the papers' text and OpenReview feedback ². This creates a valuable benchmark for future studies that want to extract and analyze limitations from research papers and study their impact on scholarly progress. Notably, due to the integration of authored text and external data, this benchmark suffers less from the self-reporting bias, providing a broader and more nuanced perspective about limitations.
- We create a novel RAG system to generate limitations. This is a useful tool for editors or researchers in the under-reporting or hedging scenario.
- Evaluating the generated limitations is a key chal-

lenge. Many traditional metrics for measuring text generation quality depend on lexical overlap (ROUGE (Lin, 2004), BLEU (Papineni, 2001)) or semantic similarity (BERTScore (Zhang et al., 2019), MoverScore (Zhao et al., 2019)). They are not sufficient for our purpose as they overprioritize common terms (e.g., tokens such as "bias", "dataset", and "generalizability" are standard in the limitations section of most AI/ML/NLP papers and should not affect the evaluation score). We develop a new evaluation framework that uses LLMs-as-judges to provide fine-grained comparisons and detailed error analysis, thus offering a more transparent and actionable evaluation for limitations.

2 Related Work

Various studies have studied the reporting, or lack thereof, of 'limitations' in scholarly papers. Ioannidis (2007) found that only 17% of top-tier journal articles mentioned limitations, with just 1% in abstracts. Similarly, Puhan et al. (Puhan et al., 2012) identified 27% of biomedical papers lacked limitations, risking overestimated research reliability. Moreover, Goodman (Goodman et al., 1994) found that in the peer review process acknowledging limitations is often problematic. Few journals require discussing limitations (Ioannidis, 2007), risking biased reviews and weakening scientific dialogue (Horton, 2002), necessitating reforms for transparency.

These gaps underscore the need for a systematic framework for limitations generation, an area that remains underexplored. Faizullah et al. (Faizullah et al., 2024) propose an LLM-chain pipeline to summarize and iteratively refine candidate limitations; Al Azher et al. (2024) integrates topic modeling with LLMs to derive and elaborate structured limitation themes; and Azher (2024) develops a graph-augmented LLM approach for generating comprehensive limitation statements. Other work focuses on overcoming limitations of charts and graph by generating meanigful captions (Al Azher and Alhoori, 2024). However, these studies are confined to ACL/EMNLP corpora and author-stated limitations, overlooking reviewer comments (e.g., OpenReview), and evaluate only with ROUGE, BERTScore, and coherence metrics, which do not assess finer-grained topic or contextual alignment. Our proposed framework addresses existing limitations by incorporating NeurIPS

¹As is the practice in the NLP and ML community in recent years.

²https://openreview.net/

papers alongside ACL, collecting OpenReview comments for each NeurIPS paper to establish a broader ground truth, and implementing a limitation-level text evaluation method to preserve granularity.

Evaluating NLP outputs is crucial for assessing quality, accuracy, and relevance. Traditional metrics like ROUGE and BLEU struggle with semantics, while BERTScore improves similarity but relies on references and lacks meaningful error analysis. Recent advancements in large language models (LLMs) have revolutionized text evaluation methods (Zheng et al., 2023) using zero-shot and in-context learning (Wei et al., 2022). For example, GPTScore (Fu et al., 2023) uses LLMs to evaluate how likely a text could be generated based on a given context, TIGERScore (Jiang et al., 2023) provide explainable error analysis, and PandaLM (Wang et al., 2023) distinguish superior models among multiple models. Other approaches, such as AttrScore (Yue et al., 2023), evaluate whether a reference supports or contradicts generated statements. SummacConv (Laban et al., 2022) breaks documents into sentence-level units and removes summary sentences with low entailment probability.

Despite their promise, LLM-based evaluations have issues such as positioning bias, where altering the order of inputs affects results. To mitigate positioning bias, we systematically alter the position of limitations and select outputs that remain consistent across different orderings. Moreover, while these text evaluation methods focus on lexical overlap and surface-level LLM-based analysis, our proposed approach advances beyond these limitations by uniquely combining granularity-aware evaluation, considering LLM as a judge, quantitative scoring, with topic-level agreement.

3 Limitation Extraction

Our limitation extraction is twofold. At first, we extracted limitations with a science parse ³ tool. If a paper has a separate limitation section, we were able to extract it using this tool. But if the limitation has no separate section, we used Python regex to check the word 'limitation' and extracted all consecutive texts before another section appeared. However, this approach also collected some noisy sentences unrelated to the limitation.

So, in the later part, we applied LLM and sent the texts to extract only limitations. We also used Selenium to collect reviews from OpenReview. Later, the 'author mentioned limitation' is sent to the zero-shot LLM to extract limitations, and the review from 'OpenReview' is sent to the zero-shot LLM to extract shortcomings/weaknesses. Through this process, we removed the noisy sentences. (§ 3.1)

3.1 Dataset of Extracted Limitations

Granularity. A key challenge in constructing a dataset of research limitations is defining the appropriate level of granularity. Should a limitation be captured as a single phrase, a full sentence, or an entire paragraph? We define a **limitation** as a *sequence of sentences*, as individual sentences often do not encapsulate multiple limitations. In contrast, a single limitation can extend across multiple sentences, sometimes forming a complete paragraph.

Extraction Sources. Two primary sources form the basis of our dataset: (1) limitations explicitly acknowledged by authors, and (2) those highlighted through peer-review commentary. Although author-reported limitations often provide well-structured insights, previous research indicates that such limitations may be underreported or carefully hedged. To address this gap, we incorporate OpenReview⁴ comments, in which peer reviewers frequently identify additional constraints or weaknesses not mentioned by the authors.

Our dataset consists of papers from conferences in natural language processing and machine learning: ACL⁵, NeurIPS⁶, and PeerJ⁷. We include 6,932 NeurIPS papers from 2021 and 2022 and 5,739 ACL papers from 2023-2024, 1000 biological and medical science papers from PeerJ. In addition, we integrate OpenReview comments for 2,802 papers from NeurIPS, thus capturing both self-reported limitations and those externally identified by reviewers. All of the PeerJ paper contains self-reported Limitations alongside other sections and peer review comments. The first experiment evaluates state-of-the-art LLMs' ability to extract limitations from papers (when they are explicitly mentioned) and their reviews. For each paper, we use GPT-40 mini to extract and get an average of 8

³https://github.com/allenai/science-parse

⁴https://openreview.net/

⁵https://aclrollingreview.org/cfp

⁶https://neurips.cc/public/guides/PaperChecklist

⁷https://peerj.com/benefits/indexing-and-impact-factor/

limitations from a paper and 10 from open review.

LLM as a Ground Truth Limitation Extractor. Identifying and isolating limitations can be challenging, as they often appear in sections such as Conclusion, Discussion, or Future Work and may also be dispersed throughout OpenReview feedback. To address this, we implement a targeted extraction process designed to be both efficient and accurate. Rather than processing entire papers, which would be computationally expensive due to larger context windows, we first apply Python regex to identify and extract relevant sections. Some research papers include a dedicated limitations section; we extract the text of this section using the AllenAI Science Parse tool⁸. When a paper doesn't have such a well-defined section, we identify sentences containing the word *limitations* and continue extracting subsequent sentences until we encounter one of the following sections: Conclusion, Future Work, Ethics, Acknowledgement, or Grant. Then, we extracted the limitations from the research paper by LLM to remove noisy sentences and make segments.

To capture broader perspectives from peer reviews, we aggregate comments from multiple OpenReview responses into a single text. Finally, we employ GPT-40 ⁹ (Figure 4) to segment the extracted text and identify distinct limitation statements, distinguishing those reported by authors from those highlighted by reviewers. We evaluate the quality of these extracted limitations through a user study (§5).

Dataset Applications. The resulting dataset will be made publicly available and can be used as a benchmark for evaluating automated limitation extraction and generation methods (§5). Beyond this, the extracted limitations can be examined and organized into a taxonomy of limitations in ML and NLP, offering a more structured understanding of common research challenges. By integrating this taxonomy into citation networks, we can introduce the concept of a Limitation Multigraph, enabling scientometric analyses into whether certain limitations shape the direction of subsequent research or, alternatively, tend to be overlooked. These avenues present new opportunities to study how the reporting (or lack thereof) of limitations affects the broader scientific discourse, a topic we plan to explore in future work.

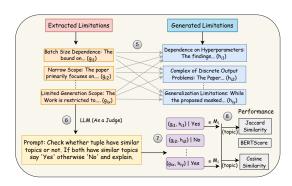


Figure 2: Evaluation of generated limitations.

4 Limitation Generation

We generated limitations by feeding every paper's text, excluding its actual "Limitations" section into both vanilla and RAG-augmented LLMs, experimenting with either the top three sections (Abstract, Introduction, Conclusion) or the full paper as input and using three ground-truth definitions (authorstated limitations, OpenReview comments, or their combination). We computed each section's cosine similarity to the "Limitations" section across ACL, NeurIPS dataset and found that Abstract, Introduction, and Conclusion rank highest on average (Table 19, Appendix). When full paper inputs exceed the context length, we make them chunkwise.

4.1 LLM and RAG for Limitation Generation

Most research papers do not explicitly mention limitations, or they underreport these shortcomings even when a dedicated section is provided. To address this gap, we propose a system that generates limitations directly from the paper text.

Vanilla LLM. In the vanilla LLM setup, the language model is prompted with the entire paper P_i . When the paper exceeds the context window, it is divided into chunks $\{P_i'\}$, and limitations are generated for each chunk (D'Arcy et al., 2024). The LLM then aggregates these chunk-specific outputs into a cohesive, meaningful final set of limitations, ensuring that no relevant material is overlooked in especially long or detailed papers.

RAG Integration. A paper P_i can be used independently to generate limitations, but this approach risks overlooking valuable insights from other, potentially related papers. In particular, even when a paper lacks an explicit limitations section, other papers with similar methodologies or datasets may discuss relevant shortcomings. More importantly, certain findings may be implicitly contradicted by subsequent research. Also, failing to consult the

⁸https://github.com/allenai/science-parse

⁹https://openai.com/index/hello-gpt-4o/

broader literature leads to a significant gap in the reported limitations. To address this issue, we employ a RAG framework, which allows the system to draw on context from multiple papers rather than relying only on P_i . Specifically, we begin by splitting 100 random papers into multiple chunks and inserting these chunks into a vector database. A retriever then identifies the top k chunks most relevant to the input text $\{P_i/P_i'\}$, depending on whether the paper P_i or a chunk P_i' can fit into the context window, using cosine similarity to match these chunks to the query (Figure 1, step 4,5).

5 Evaluation of Generated Limitations

We introduce a framework, called *PointWise* (**PW**), for evaluating the performance of LLMs in generating limitations for scholarly papers (Figure 2).

Problem Setup. Suppose we have a set of papers $P = \{P_1, P_2, \dots, P_n\}$. For each paper P_i , we assume access to: Ground truth limitations $G_i = \{g_{i1}, g_{i2}, \dots, g_{ix}\}$, where x is the number of ground truth limitations we extracted or annotated for P_i . And LLM-generated limitations $H_i = \{h_{i1}, h_{i2}, \dots, h_{iy}\}$, where y is the number of limitations produced by the LLM for P_i . Our goal is to measure (1) how many ground truth limitations the LLM correctly reproduces (coverage) and (2) how well each matched pair of limitations aligns in content and focus (performance).

5.1 Coverage

A. Pairwise Matching. To quantify coverage, we first create all possible pairs of limitations between the sets G_i and H_i . Let

$$S_i = \{(g_{ik}, h_{il}) \mid 1 \le k \le x, 1 \le l \le y\}.$$

Hence, $|S_i| = x \times y$. We then use an LLM as a judge (Zheng et al., 2023) to decide if a ground truth limitation g_{ik} and a generated limitation h_{il} are similar in content or topic:

$$J(g_{ik}, h_{il}) = \begin{cases} 1, & \text{if } g_{ik} \text{ and } h_{il} \text{ are similar,} \\ 0, & \text{otherwise.} \end{cases}$$

We collect all *matched* pairs into a set

$$M_i = \{(g_{ik}, h_{il}) \mid J(g_{ik}, h_{il}) = 1\},\$$

and let $|M_i| = z_i$ be the number of matched pairs for paper P_i .

B. Coverage of Ground Truth Limitations. We define $C_{Gi}(g_{ik}) = 1$ if the ground truth limitation g_{ik} appears in *at least one* matched pair in M_i , and 0 otherwise:

$$C_{Gi}(g_{ik}) = \begin{cases} 1, & \exists h_{il} \text{ such that } (g_{ik}, h_{il}) \in M_i, \\ 0, & \text{otherwise.} \end{cases}$$

The coverage of ground truth limitations for paper P_i is

$$A_{Gi} = \frac{1}{x} \sum_{k=1}^{x} C_{Gi}(g_{ik}).$$

In other words, A_{Gi} measures the fraction of ground truth limitations in P_i that are matched with at least one LLM-generated limitation.

C. Coverage of LLM-Generated Limitations. Similarly, we define $C_{Hi}(h_{il}) = 1$ if a generated limitation h_{il} appears in *at least one* matched pair in M_i , and 0 otherwise:

$$C_{Hi}(h_{il}) = \begin{cases} 1, & \exists g_{ik} \text{ such that } (g_{ik}, h_{il}) \in M_i, \\ 0, & \text{otherwise.} \end{cases}$$

The coverage of LLM-generated limitations for paper P_i is

$$A_{Hi} = \frac{1}{y} \sum_{l=1}^{y} C_{Hi}(h_{il}).$$

We aggregate these coverage values across all papers by taking their means:

$$A_G = \frac{1}{n} \sum_{i=1}^n A_{Gi}, \quad A_H = \frac{1}{n} \sum_{i=1}^n A_{Hi}.$$

D. Precision, Recall, and F₁. We also compute overall precision, recall, and F₁ scores. For each paper P_i :

$$TP_i = |M_i|, FP_i = x - \sum_{k=1}^{x} C_{Gi}(g_{ik}),$$

$$FN_i = y - \sum_{l=1}^{g} C_{Hi}(h_{il}).$$

Here, TP_i (true positives) is the total number of matched pairs; FP_i (false positives) is the number of ground truth limitations not matched by any LLM-generated limitation; FN_i (false negatives) is the number of LLM-generated limitations unmatched by any ground truth limitation. True negatives (TN_i) are not defined since no "negative"

class of limitations is specified. If there is one ground truth limitation g_{i1} matches with multiple LLM-generated limitations (and vice versa), True Positive (TP) counts as one. (Details in Appendix A.1)

5.2 Performance

After finding matched pairs $(g_{ik}, h_{il}) \in M_i$, we evaluate their *quality* using text based evaluation (Rouge-1, Rouge-2, BERTSCore, Cosine Similarity), keyword based evaluation, and heading based evaluation (Details in Appendix A.2).

6 Experimental Setup

We fine-tuned three sequence-to-sequence models, T5 (512-token window) (Raffel et al., 2020), BART (1024 tokens), and Pegasus (1024 tokens) (Zhang et al., 2020) on a 70 / 30 train-test split. All were trained for 3 epochs with a learning rate of 5×10^{-5} , weight decay 0.01, 300 warmup steps, and batch sizes of 4 (train) and 8 (eval), with early stopping; inputs longer than 512 tokens were truncated. For zero-shot evaluation, we used, GPT-3.5 (16 K) and GPT-40 (128 K) at temperature 0.1, and Llama 3 8B and Mistral 7B (128 K) via Ollama. To incorporate retrieval (RAG), we built a VectorStoreIndex with llama-index 10, and OpenAI text-embedding-ada-002 embedding model for encoding the source and query documents. We took 100 random papers from the dataset and stored them in a vector database. These documents were chunked into 800-token segments with 20-token overlaps to preserve context, and the LLM temperature was set to 0.

7 Experiments and Results

7.1 Limitation Extraction Evaluation

We manually collected limitation sections from 500 ACL papers randomly (human curated). We measure the n-gram overlap and semantic similarity between human-curated limitations vs. automatically extracted ground truth limitations to evaluate our automatic limitation extraction process. As shown in Table 1, our automatically extracted limitation text achieves strong alignment with human-curated ground truth limitations on the ACL '23 and ACL '24 datasets, with ROUGE-1, L scores exceeding 85% and BERTScore above 95%.

Metrics	ACL 24	ACL 23
ROUGE-1	86.68	85.07
ROUGE-2	83.13	81.68
ROUGE-L	86.64	85.04
BLEU	77.70	75.15
Jaccard Similarity	79.36	77.08
JS Divergence	23.58	24.88
Levenshtein	92.54	90.32
Cosine Similarity	90.05	88.70
BERTScore (F1)	95.70	95.66

Table 1: Evaluating the Alignment of Human-Collected vs. Automated Limitation Extraction.

Human Evaluation. We employ three annotators (separate from this paper's authors)¹¹ to evaluate the quality of the extracted limitations. We choose a sample of 100 limitations, and for each, we show them the source and ask a Yes/No question, whether they thought the LLM extracted the limitation from the source, as opposed to generating it from scratch. Each annotator answer positively in > 90% of cases (avg \pm std= $95 \pm 2.45\%$) (Table 2). The annotator pairs have weighted Kappa scores between 38 - 75%, denoting fair to substantial agreement (Table 18, Appendix). We checked Llama 3 70b as an extractor, but the extraction quality was not good. Overall, the extraction quality is high, and the annotators agree; therefore, we use the extracted limitations as ground truth in the subsequent evaluations. Furthermore, two annotators manually verified the extracted limitations from 1000 papers. They assessed whether each automatically extracted limitation (via tool + LLM) was grounded in the original paper text or OpenReview content and whether it included any irrelevant or unrelated information. Their analysis confirmed that all extracted limitations were faithfully sourced, with no instances of hallucinated, noisy, or newly generated content.

Model	Role	Sample	U1	U2	U3
GPT 40 mini	Extractor	100	92	95	98

Table 2: Evaluating LLM as an extractor role with human annotator (U).

7.2 Limitation Generation Evaluation

Coverage Evaluation. Table 3 reports results on ACL data using only author-stated limitations as ground truth. We couldn't collect the Open-Review of ACL papers, so we used the authormentioned limitations as ground truth. GPT-3.5

¹⁰https://www.llamaindex.ai/.

 $^{^{11}\}mbox{CS}$ graduate students with research experience in NLP and AI

R-L	BS	JS_{text}	CS _{text}	C_{GT}	C_{LLM}	Prec.	Recall	F1
19.92	87.81	10.82	31.79	35.48	29.59	0.29	0.31	0.30
19.43	87.67	10.68	31.91	33.71	30.10	0.30	0.31	0.31
20.15	87.66	10.71	33.39	29.28	25.27	0.25	0.26	0.26
25.66	88.30	14.69	40.4	61.38	39.04	0.39	0.50	0.44
24.24	87.08	14.65	43.12	76.62	46.65	0.47	0.67	0.55
30.21	90.88	19.47	45.37	39.99	44.66	0.42	0.40	0.41
16.57	86.02	8.70	32.29	57.65	19.76	0.20	0.31	0.24
23.17	87.33	12.99	39.29	67.13	45.67	0.57	0.45	0.51
	19.92 19.43 20.15 25.66 24.24 30.21 16.57	19.92 87.81 19.43 87.67 20.15 87.66 25.66 88.30 24.24 87.08 30.21 90.88 16.57 86.02	19.92 87.81 10.82 19.43 87.67 10.68 20.15 87.66 10.71 25.66 88.30 14.69 24.24 87.08 14.65 30.21 90.88 19.47 16.57 86.02 8.70	19.92 87.81 10.82 31.79 19.43 87.67 10.68 31.91 20.15 87.66 10.71 33.39 25.66 88.30 14.69 40.4 24.24 87.08 14.65 43.12 30.21 90.88 19.47 45.37 16.57 86.02 8.70 32.29	19.92 87.81 10.82 31.79 35.48 19.43 87.67 10.68 31.91 33.71 20.15 87.66 10.71 33.39 29.28 25.66 88.30 14.69 40.4 61.38 24.24 87.08 14.65 43.12 76.62 30.21 90.88 19.47 45.37 39.99 16.57 86.02 8.70 32.29 57.65	19.92 87.81 10.82 31.79 35.48 29.59 19.43 87.67 10.68 31.91 33.71 30.10 20.15 87.66 10.71 33.39 29.28 25.27 25.66 88.30 14.69 40.4 61.38 39.04 24.24 87.08 14.65 43.12 76.62 46.65 30.21 90.88 19.47 45.37 39.99 44.66 16.57 86.02 8.70 32.29 57.65 19.76	19.92 87.81 10.82 31.79 35.48 29.59 0.29 19.43 87.67 10.68 31.91 33.71 30.10 0.30 20.15 87.66 10.71 33.39 29.28 25.27 0.25 25.66 88.30 14.69 40.4 61.38 39.04 0.39 24.24 87.08 14.65 43.12 76.62 46.65 0.47 30.21 90.88 19.47 45.37 39.99 44.66 0.42 16.57 86.02 8.70 32.29 57.65 19.76 0.20	19.92 87.81 10.82 31.79 35.48 29.59 0.29 0.31 19.43 87.67 10.68 31.91 33.71 30.10 0.30 0.31 20.15 87.66 10.71 33.39 29.28 25.27 0.25 0.26 25.66 88.30 14.69 40.4 61.38 39.04 0.39 0.50 24.24 87.08 14.65 43.12 76.62 46.65 0.47 0.67 30.21 90.88 19.47 45.37 39.99 44.66 0.42 0.40 16.57 86.02 8.70 32.29 57.65 19.76 0.20 0.31

Table 3: Results of models in "Coverage" (Coverage of Ground Truth Limitation (C_{GT}), LLM Generated Limitation (C_{LLM}), Precision, Recall, and F1-score) and "performance" metrics – **R**ouge-x, BLEU, BertScore (BS), Jaccard (JS) and Cosine (CS) similarity over text (ACL data).

achieves the highest Coverage of Ground Truth (C_{GT} =76.62) and Coverage of LLM-Generated Limitations (C_{LLM} =46.65). Vanilla LLMs tend to produce overly verbose limitations. Integrating RAG yields more focused, concise outputs that cover a broader range of topics and align better with ground truth—for both GPT-3.5 and GPT-40 mini. In particular, incorporating RAG increases C_{LLM} by 25.91 points over GPT-40 mini. For NeurIPS data (Table 4), considering RAG increase C_{GT} , decrease C_{LLM} , which gains F1 score (+0.06). Compared to ACL, NeurIPS, and Peerj, the ACL dataset archives the highest C_{GT} , and NeurIPS achieves the highest C_{LLM} (Table 11, Appendix).

Performance Evaluation. Table 3 reports our performance metrics using author-extracted limitations as ground truth. Compared to all models, GPT 3.5 + RAG shows the best results. Augmenting GPT-3.5 with RAG boosts every measure, ROUGE-L (+5.97) and BERTScore (+3.8). GPT-40 mini exhibits the same upward trend when paired with RAG. For NeurIPS data (Table 4), considering RAG increased the performance of all metrics except Cosine Similarity. Among our datasets, NeurIPS achieves the highest ROUGE, BERTScore, and BLEU largely due to the inclusion of OpenReview feedback (Table 11, Appendix). Moreover, our proposed PointWise approach achieves better BERTScore than the traditional way in every model (Table 8, 17, Appendix).

Evaluation 'LLM as a Judge'. We evaluated our Judge LLM model GPT 40 mini regarding PointWise evaluation approach. From each pair of ground truth and LLM generated limitation, we asked LLM whether both are same or not (Figure 3). We use the same annotators as before and choose a sample of 200 pairs, 100 *positive* cases where the LLM predicted a match and 100 *negatives*. We ask humans the same question as the model and calculate the weighted Kappa scores.

Metric	GPT 40 mini	GPT 40 mini + RAG
Rouge-L	17.07	18.10 (+1.03)
BERTScore	86.25	86.45 (+0.2)
JS_{text}	9.12	9.72 (+0.6)
CS_{text}	35.42	34.71 (-0.71)
C_{GT}	45.69	61.76 (+16.07)
C_{LLM}	76.27	59.18 (-17.09)
F1 score	0.54	0.60 (+0.06)

Table 4: Coverage and Performance measure using GPT 40 mini + RAG in NeurIPS data. Ground truth is Limitation + OpenReview, and input is the top 3 sections alongside the abstract.

As seen in (Table 9, Appendix), the human experts overwhelmingly agree with GPT, with kappa values ranging between 90-95%. The agreement between human annotators is > 92%, indicating a perfect agreement. We also experimented with a strong baseline of Llama-3.1 400B, which shows a poor agreement with humans.

7.3 Ablation Study

We ablate in two dimensions a) the size of the input text. b) the ground truth

a. Size of the input text: Table 5 presents results using GPT-40 mini. Moving from the top three sections to all sections as input improves every metric across all ground-truth types. The only exceptions are a few cases, such as BERTScore and LLM Coverage (C_{LLM}) when using OpenReview as a ground truth. We also experimented ablation study with other models (Table 6), such as Mistral 7B and Llama 3 8B on NeurIPS data. Using all sections as input yields more detailed outputs, boosting LLM coverage by 8.35 points. However, it slightly reduces ground-truth coverage, C_{GT} (-0.29), ROUGE-L (+0.64), and BERTScore (+0.36). Llama 3 8B shows the opposite trend: using all sections significantly degrades its performance. This likely reflects Llama's limited context window, which hinders its ability to process longer inputs and generate coherent limitations.

b. Ground truth: Table 5 compares zero-shot GPT-40 mini results against three ground-truth sets: author-stated limitations, OpenReview comments, and their combination. Using both Limitations + OpenReview yields higher scores on nearly every metric compared to Limitations alone, except for CS_{Keyword}, JS_{Keyword}, and C_{GT}. This demonstrates that reviewer feedback enriches the ground truth, improving overall performance, but because the combined ground truth is larger, its relative coverage (C_{GT}) can decrease. We also experimented with Mistral 7b (Table 7, Appendix), incorporating OpenReivew increased C_{LLM}, Cosine similarity, and F1 score but marginal drops in other metrics.

Metric	Input		Ground Tr	ruth
1,100110	Sec	Auth	OPR	Auth + OPR
R-L	3	16.57	18.55	17.07 (+0.5)
R-L	All	17.11 (1)	18.65 (₁)	18.83 (+1.72)
BS	3	86.02	86.50	86.25 (+0.23)
BS	All	86.24 (1)	85.96 (1)	86.67 (+0.43)
JS_{text}	3	8.70	10.17	9.12 (+0.42)
JS_{text}	All	9.22 (1)	10.23 (₁)	10.38 (+1.16)
CS_{text}	3	32.29	38.35	35.42 (+3.13)
CS_{text}	All	33.23 (↓)	38.74 ()	38.87 (+5.64)
C_{GT}	3	57.65	54.79	45.69 (-12.41)
C_{GT}	All	61.55 (_†)	56.55 (1)	49.43 (-12.12)
C_{LLM}	3	19.76	56.79	76.27 (+57.07)
CLLM	All	40.04 (1)	56.77 (1)	42.59 (+2.55)
F1	3	0.24	0.55	0.54 (+0.3)
F1	All	0.45 (1)	0.56 (↑)	0.44 (-0.01)

Table 5: GPT 40 mini results in "coverage" (Coverage of Ground Truth Limitation (C_{GT}), LLM Generated Limitation (C_{LLM}), F1-score) and "performance" (**R**ouge-x, **B**ert**S**core (over entire text and just heading), Jaccard (JS) and Cosine (CS) similarity over text and heading) (NeurIPS data).

8 Discussion

Our results demonstrate that leveraging external OpenReview feedback alongside author-stated limitations substantially enriches the scope and quality of LLM-generated limitations, improving performance and coverage. This indicates the model not only reproduces author-reported issues but also uncovers additional shortcomings highlighted by peer reviewers. Similarly, integrating a Retrieval-Augmented Generation (RAG) module with GPT-3.5/4 improved every metric by grounding the LLM in relevant external context. The RAG-enhanced model produced more concise, targeted limitations (fewer sentences on average), trading a slight dip in recall with improved performance. Using all sections marginally improved performance then considering 3 section in most of the metrics. Moreover,

Metric	Input (3 Sec.)	Input (All Sec.)
Mistral 7B		
Rouge-L	14.59	13.95 (-0.64)
BERTScore	84.59	84.23 (-0.36)
JS_{text}	7.70	7.49 (-0.21)
CS_{text}	30.95	31.82 (+0.87)
BLEU	0.30	0.30(0)
C_{GT}	38.28	37.99 (-0.29)
C_{LLM}	22.09	30.44 (+8.35)
F1 score	0.24	0.32 (+0.08)
Llama 3 8B		
Rouge-L	25.66	22.78 (-2.88)
BERTScore	88.30	87.29 (-1.01)
JS_{text}	14.69	11.16 (-3.53)
CS_{text}	40.40	34.80 (-5.60)
$JS_{keyword}$	20.28	20.31 (+0.03)
C_{GT}	61.38	18.01 (-43.37)
C_{LLM}	39.04	5.01 (-34.03)
F1 score	0.44	0.05 (-0.39)

Table 6: Ablation study comparing Mistral 7B and Llama 3 8B on NeurIPS data. We report metrics with only the top 3 sections versus all sections as input, using author-mentioned limitations as ground truth.

our PointWise (PW) framework breaks each limitation into fine-grained points. It employs an LLM "judge" to evaluate alignment one by one, where traditional lexical metrics can't capture depth and relevance. Finally, chunking long documents ensures no critical content is lost outside an LLM's context window, and RAG helps focus on significant, relevant limitations outside the current paper. Moreover, models like GPT-40 mini deliver consistently strong performance with high correlation with humans as extractors and judges, whereas others (e.g., Llama) occasionally introduce noise. These findings affirm that combining external review data, targeted section selection, RAG augmentation, and a detailed PointWise evaluation yields more trustworthy and comprehensive LLMgenerated limitations—and paves the way for future work in automated scientific limitations.

9 Conclusion

We present a new approach for automatically extracting, generating and evaluating limitations in scientific articles. Our method explores chunk analysis, accommodating top sections of the entire paper and integrating OpenReview feedback to capture perspectives beyond those of the original authors. To evaluate the effectiveness of our system, we introduce a granular text evaluation framework that breaks down limitations into more minor points and employs LLMs as a Judge for assessing topic-

level alignment. Combining standard NLP metrics with a keyword-based analysis provides a more nuanced and reliable assessment. Our findings are further validated by human review, which reveals a strong alignment between human judgments and the LLM-driven extraction process, including when the LLM is used as a judge.

Limitations

In this work, we focused on venues in natural language processing (ACL papers from 2023-2024) and machine learning (NeurIPS papers 2021-2022), and Biology domain papers from PeerJ, which ensures high relevance and quality but insufficient for broader generalizability. While this scope allows us to benchmark the performance of LLMs in extracting limitations from well-structured scientific texts, we acknowledge that the findings may not generalize to papers from other fields, such as social sciences, physics, chemistry, or mathematics where writing conventions and limitation styles may differ. Due to high API costs, we did not experiment with GPT-4 or GPT-4o; instead, we opted for GPT-40 Mini as a cost-effective alternative. Similarly, to avoid the high expense of maintaining an extensive vector database for RAG with closed LLMs, we limited our retrieval corpus to 100 randomly selected papers. We did not evaluate Open LLMs with RAG, such as Llama and Mistral, to assess their performance in retrieval-augmented settings. While we incorporated OpenReview comments for NeurIPS papers, we could not find them for ACL papers. Furthermore, we relied only on GPT-40 Mini as the evaluation judge and did not experiment with other LLMs for assessment. To evaluate the effectiveness of LLMs as both text extractors and judges, we conducted a human annotation study with 200 samples and only three annotators.

We used GPT 40 mini for text extraction, text generation, and LLM as a judge. We experimented with other LLMs such as Llama but GPT 40 mini performs better. To rely on LLM for limitation extraction and evaluation ther are some potential biases.

- Self-Validation Bias: We used the same LLM for text generation and evaluation in experiments. Where a model can favor its style or logic.
- Positional bias: For LLM as a judge, we have

- multiple texts to check whether both are similar. Here, LLM can have positional bias, which is more centered on the first input text.
- Confirmation bias: We experimented with the same LLM (GPT 40 mini) for text extraction and generation. Even though both are used in an isolated environment with zero shot (no fine-tuning), the same LLM can produce confirmation bias on some common/generic limitations such as "Generalizability issue" or "small sample size" which reduces novelty.
- Hallucination and Overconfidence: When extracting limitations, LLM might sometimes invent limitations that aren't grounded in the source text.

To mitigate self-validation bias, we compared our LLM's judgments against human evaluation, tested alternative models for generation and judgment (Llama and Mistral), and incorporated RAG. We changed the order and considered the consistent one to overcome the positional bias problem. Moreover, we integrated RAG with LLM for text generation, where LLM can get diverse evidence or force diversity sampling on generations to minimize the confirmation bias problem. To overcome hallucination and overconfidence, we evaluate the extraction process by LLM with three annotators to determine whether LLM extracts limitations without inventing anything new. However, these approaches are not sufficient.

As part of our future work, we plan to expand the dataset to include papers from a broader range of domains, such as arXiv submissions from diverse subject areas (e.g., bioinformatics, cognitive science) and other peer-reviewed venues beyond ACL and NeurIPS. We also aim to incorporate more review texts from platforms, where available, across different conferences. This expansion will allow us to test the robustness of our limitation extraction models across domains and improve generalizability. We will explore broader contexts from open LLMs incorporating RAG, and integrating more extensive multi-agent LLM frameworks. Moreover, we will gather annotations from a larger and more diverse pool of experts, enabling a deeper analysis of generated limitations. In doing so, we aim to make our framework more robust and broadly applicable while maintaining its core strength of providing granular, interpretable assessments of LLM-generated limitations.

Ethics Statement

This research adheres to the ethical standards in scientific research. All data used in this study, including research papers and OpenReview feedback, were obtained from publicly available sources and used in compliance with their respective usage policies. No private or sensitive information was accessed or utilized. We acknowledge the potential biases inherent in large language models (LLMs) and have taken steps to mitigate their impact. For instance, we conducted human evaluations to validate the outputs of LLMs, ensuring alignment with human judgment and reducing the risk of generating misleading or inaccurate content. Additionally, our evaluation framework emphasizes transparency and granularity, enabling a more nuanced assessment of generated limitations. For the user study, we engaged three computer science graduate students; they are no longer involved in this research. Their participation was voluntary, and they had no conflicts of interest regarding the study. We are committed to promoting the responsible use of AI in scientific research. Our work aims to enhance scientific studies' transparency, reliability, and reproducibility by generating and evaluating limitations systematically and ethically. We encourage further research to address potential biases in LLMs and improve their application in scientific contexts. Our study didn't filter out any dataset based on race, ethnicity, color, location, gender, or other discriminatory attributes. Therefore, we are dedicated to complying with the ethical code and policy attributed to ACL. Large language models inherently carry biases and risks derived from their training corpora—confirmation bias toward "safe" limitations (e.g., small sample size), fluency bias favoring well-written over factually grounded outputs, verbosity bias rewarding longer responses, and selfvalidation bias when the same model is used for extraction, generation, and evaluation. To mitigate these, we ground all generations in source content and peer reviews via Retrieval-Augmented Generation (RAG), diversify our ground truth by incorporating human-authored OpenReview critiques, and conduct parallel human evaluations to detect overconfidence and evaluation bias. RAG-based outputs are demonstrably more concise, reducing verbosity bias, and our multi-model judge setup (including Llama and Mistral alongside GPT-40 mini) helps break circularity. We acknowledge that further work is needed to quantify these biases

rigorously and plan to investigate cross-domain robustness and fairness in future studies.

References

- Ibrahim Al Azher and Hamed Alhoori. 2024. Mitigating visual limitations of research papers. In 2024 *IEEE International Conference on Big Data (Big-Data)*, pages 8614–8616. IEEE.
- Ibrahim Al Azher, Venkata Devesh Reddy, Hamed Alhoori, and Akhil Pandey Akella. 2024. Limtopic: Llm-based topic modeling and text summarization for analyzing scientific articles limitations. In ACM/IEE Joint Conference on Digital Libraries (JCDL).
- Thomas M Annesley. 2010. The discussion section: your closing argument. *Clinical chemistry*, 56(11):1671–1674.
- Ibrahim Al Azher. 2024. Generating suggestive limitations from research articles using llm and graph-based approach. In *Proceedings of the 24th ACM/IEEE Joint Conference on Digital Libraries*, pages 1–3.
- Ibrahim Al Azher, Miftahul Jannat Mokarrama, Zhishuai Guo, Sagnik Ray Choudhury, and Hamed Alhoori. 2025. Futuregen: Llm-rag approach to generate the future work of scientific article. *arXiv* preprint arXiv:2503.16561.
- Suzanne Bunniss and Diane R Kelly. 2010. Research paradigms in medical education research. *Medical education*, 44(4):358–366.
- Lisa Chasan-Taber. 2014. Writing dissertation and grant proposals: Epidemiology, preventive medicine and biostatistics. CRC Press.
- Mike D'Arcy, Tom Hope, Larry Birnbaum, and Doug Downey. 2024. Marg: Multi-agent review generation for scientific papers. *arXiv preprint arXiv:2401.04259*.
- Jean Escande, Christophe Proust, and Jean Christophe Le Coze. 2016. Limitations of current risk assessment methods to foresee emerging risks: Towards a new methodology? *Journal of Loss Prevention in the Process Industries*, 43:730–735.
- Kevin W Eva and Lorelei Lingard. 2008. What's next? a guiding question for educators engaged in educational research.
- Abdur Rahman Bin Mohammed Faizullah, Ashok Urlana, and Rahul Mishra. 2024. Limgen: Probing the llms for generating suggestive limitations of research papers. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 106–124. Springer.
- Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. Gptscore: Evaluate as you desire. *arXiv* preprint arXiv:2302.04166.

- Steven N Goodman, Jesse Berlin, Suzanne W Fletcher, and Robert H Fletcher. 1994. Manuscript quality before and after peer review and editing at annals of internal medicine. *Annals of internal medicine*, 121(1):11–21.
- Sue Greener. 2018. Research limitations: the need for honesty and common sense.
- Maarten Grootendorst. 2020. Keybert: Minimal keyword extraction with bert.
- Richard Horton. 2002. The hidden research paper. *Jama*, 287(21):2775–2778.
- Ken Hyland. 1998. Hedging in scientific research articles.
- John PA Ioannidis. 2007. Limitations are not properly acknowledged in the scientific literature. *Journal of clinical epidemiology*, 60(4):324–329.
- Dongfu Jiang, Yishan Li, Ge Zhang, Wenhao Huang, Bill Yuchen Lin, and Wenhu Chen. 2023. Tigerscore: Towards building explainable metric for all text generation tasks. *Transactions on Machine Learning Research*.
- Philippe Laban, Tobias Schnabel, Paul N Bennett, and Marti A Hearst. 2022. Summac: Re-visiting nlibased models for inconsistency detection in summarization. *Transactions of the Association for Computational Linguistics*, 10:163–177.
- Peder Larsen and Markus Von Ins. 2010. The rate of growth in scientific publication and the decline in coverage provided by science citation index. *Scientometrics*, 84(3):575–603.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Victor M Montori, Roman Jaeschke, Holger J Schünemann, Mohit Bhandari, Jan L Brozek, PJ Devereaux, and Gordon H Guyatt. 2004. Users' guide to detecting misleading claims in clinical research reports. *Bmj*, 329(7474):1093–1096.
- Kishore Papineni. 2001. Bleu: a method for automatic evaluation of mt. *Research Report, Computer Science RC22176 (W0109-022)*.
- MA Puhan, N Heller, I Joleska, L Siebeling, P Muggensturm, M Umbehr, S Goodman, and G ter Riet. 2009. Acknowledging limitations in biomedical studies: The alibi study. In *The Sixth International Congress on Peer Review and Biomedical Publication*, pages 10–12. JAMA and BMJ Vancouver, Canada.
- Milo A Puhan, Elie A Akl, Dianne Bryant, Feng Xie, Giovanni Apolone, and Gerben ter Riet. 2012. Discussing study limitations in reports of biomedical studies-the need for more transparency. *Health and quality of life outcomes*, 10:1–4.

- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67.
- Paula T Ross and Nikki L Bibler Zaidi. 2019. Limited by our limitations. *Perspectives on medical education*, 8:261–264.
- Gerben Ter Riet, Paula Chesley, Alan G Gross, Lara Siebeling, Patrick Muggensturm, Nadine Heller, Martin Umbehr, Daniela Vollenweider, Tsung Yu, Elie A Akl, et al. 2013. All that glitters isn't gold: a survey on acknowledgment of limitations in biomedical studies. *PloS one*, 8(11):e73623.
- Yidong Wang, Zhuohao Yu, Zhengran Zeng, Linyi Yang, Cunxiang Wang, Hao Chen, Chaoya Jiang, Rui Xie, Jindong Wang, Xing Xie, et al. 2023. Pandalm: An automatic evaluation benchmark for llm instruction tuning optimization. *arXiv preprint arXiv:2306.05087*.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*.
- Xiang Yue, Boshi Wang, Ziru Chen, Kai Zhang, Yu Su, and Huan Sun. 2023. Automatic evaluation of attribution by large language models. *arXiv* preprint *arXiv*:2305.06311.
- Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter Liu. 2020. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization. In *International conference on machine learning*, pages 11328–11339. PMLR.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.
- Wei Zhao, Maxime Peyrard, Fei Liu, Jing Gao, Christian M. Meyer, and Steffen Eger. 2019. Moverscore: Text generation evaluating with contextualized embeddings and earth mover distance. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 563–578.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623.
- Vilma Žydžiūnaitė. 2018. Implementing ethical principles in social research: Challenges, possibilities and limitations. *Profesinis rengimas: tyrimai ir realijos*, 1(29):19–43.

A Appendix

In our PointWise evaluation method, we measured precision, recall, and F1 score from True Positive, False Positive, and False Negative.

A.1 Coverage Measurement

We compute:

$$\mathbf{P}_{r_i} \; = \; rac{\mathrm{TP}_i}{\mathrm{TP}_i + \mathrm{FP}_i}, \quad \mathbf{R}_{r_i} \; = \; rac{\mathrm{TP}_i}{\mathrm{TP}_i + \mathrm{FN}_i},$$

and the F_1 score is the harmonic mean of P_{r_i} and R_{r_i} .

A.2 Performance Measurement

A. Text-Based Evaluation. We apply standard text similarity metrics to each matched pair, including ROUGE-1, ROUGE-L, BERTScore, Cosine Similarity, Jaccard Similarity, and BLEU, calculating the number of overlapping unigrams, the longest sequence of words, and the similarity between contextual embeddings.

B. Keyword-Based Evaluation. We employ KeyBERT (Grootendorst, 2020) to extract a set of top keywords from the ground truth limitations K_{G_i} and from the LLM-generated limitations K_{H_i} . We then measure the cosine and Jaccard similarity between K_{G_i} and K_{H_i} for each paper P_i and average these scores across the dataset.

C. Heading-Based Evaluation. We also compare concise "headings" or short titles for each limitation. Let T_{G_i} be the heading for G_i and T_{H_i} the heading for H_i . We compute BERTScore between T_{G_i} and T_{H_i} for every paper P_i and then average these values. This provides a high-level measure of how closely the top-level concepts align.

By combining coverage and performance metrics in a PointWise manner, our framework provides a detailed assessment of how well an LLM-generated set of limitations captures the breadth and depth of the ground truth. This approach also facilitates fine-grained error analysis by examining matched pairs on a per-limitation basis.

We measure coverage for both ground truth and LLM-generated limitations *independently*, focusing on each unique limitation within the matched pairs.

Furthermore, we conduct experiments using:

1. The top three sections (*Abstract, Introduction, and Conclusion*)

2. The entire paper (full paper)

This setup enables us to examine how restricting the analysis to specific sections affects coverage and matching performance.

We used three distinct prompts to check the topic level similarity between ground truth limitations and LLM-generated limitations (Figure 3, Appendix). To overcome the position bias, we choose the consistent one.

We also compared BERTScore between the Traditional approach vs the Pointwise approach, where the Pointwise approach preserves nuanced details and showed better performance (Table 8, Appendix).

Prompt 1 = '''

A tuple contains (list1, list2). Check whether both 'list1' and 'list2' have similar topics or limitation. If both have similar topics or limitations you can say "Yes", otherwise "No". Your answer should be "Yes" or "No" with explanation.

Prompt 2 = '''

A tuple contains (list2, list1). Check whether both 'list2' and 'list1' have similar topics or limitation. If both have similar topics or limitations you can say "Yes", otherwise "No". Your answer should be "Yes" or "No" with explanation.

Prompt 3 = '''

Check whether 'list2' contains a topic or limitation from 'list1' or 'list1' contains a topic or limitation from 'list2'. Your answer should be "Yes" or "No" with explanation.

Figure 3: LLM as a Judge for each limitation. Note: We use three distinct prompts to verify consistency.

Metric	Limitation	Limitaiton + OPR
Rouge-L	14.59	14.52 (-0.07)
BERTScore	84.59	84.54 (-0.05)
JS_{text}	7.70	7.62 (-0.08)
CS_{text}	30.95	32.02 (+1.07)
C_{GT}	38.28	32.51 (-5.77)
C_{LLM}	22.09	44.46 (+22.37)
F1 score	0.24	0.34 (+0.1)

Table 7: Ablation study with Mistral 7b in NeurIPS data, considering Limitation and Limitation + OpenReview (OPR) as ground truth. Top 3 sections are input here.

Approach	Traditional way	Pointwise approach
T5	82.60	87.81 (+5.21)
BART	83.09	87.67 (+ 4.58)
Pegasus	82.76	87.66 (+4.9)
GPT 3.5	83.79	87.08 (+3.29)
GPT 40 mini	82.66	86.02 (+3.36)

Table 8: Comparison of Traditional Approach and Pointwise Approach. BERTScore of each model in ACL data (Limitation is ground truth)

	GPT-4	Llama	HE1	HE2	HE3
GPT-4	1	0.71	0.9	0.92	0.95
Llama	-	1	0.81	0.79	0.76
HE1	-	-	1	0.98	0.95
HE2	-	-	-	1	0.97
HE3	-	-	-	-	1

Table 9: Evaluating how good LLM 'as a Judge' by checking Human Expert (HE) and model (GPT-40 mini, Llama-3.1 400B) agreement in determining whether an extracted limitation *matches* a generated one (in PointWise Evaluation).

Model	R-1	R-2	R-L	BertScore	BertScore _{heading}	JS _{text}	CS _{text}	BLEU	CSkeyword	JSkeyword
T5	27.13	5.11	19.92	87.81	-	10.82	31.79	0.4	18.69	11.26
BART	26.05	5.49	19.43	87.67	-	10.68	31.91	0.7	19.89	12.19
Pegasus	27.0	5.46	20.15	87.66	-	10.71	33.39	0.6	17.65	10.58
GPT 3.5	32.74	10.92	24.24	87.08	85.50	14.65	43.12	3.54	26.54	17.46
Llama 3	33.98	12.02	25.66	88.30	85.30	14.69	40.4	3.6	30.33	20.28
GPT 4o-mini	22.90	3.99	16.57	86.02	75.02	8.70	32.29	0.33	16.20	9.85
GPT 3.5 + RAG	38.11	16.31	30.21	90.88	87.65	19.47	45.37	6.64	35.74	24.65
GPT 40 mini + RAG	31	9.95	23.17	87.33	-	12.99	39.29	2.69	26.47	17.19

Table 10: Results of models in "performance" metrics – **R**ouge-x, BLEU, BertScore (over entire text and just heading), Jaccard (JS) and Cosine (CS) similarity over text and heading (ACL data).

Metric	ACL	NeurIPS	PeerJ
Rouge-L	16.57	17.07	10.62
BERTScore	86.02	86.25	83.26
JS_{text}	8.70	9.12	5.81
CS_{text}	32.29	35.42	26.19
$CS_{keyword}$	16.20	15.67	2.89
$JS_{keyword}$	9.85	9.31	1.66
BLĖU	0.33	0.4	0.01
C_{GT}	57.65	45.69	55.07
C_{LLM}	19.76	76.27	35.87
F1 score	0.24	0.54	0.39

Table 11: Comparison among ACL and NeurIPS PeerJ data in terms of Coverage and Performance. Full text is input, using GPT 40 mini to generate limitations, and Ground truth is Limitation + Peer Review (for NeurIPS and PeerJ), ACL has only author-mentioned Limitations.

Metric	Input		Ground Tr	uth
1120110	Sec	Auth	OPR	Auth + OPR
R-1	3	22.90	26.13	23.69 (+0.79)
R-1	All	23.74 (1)	26.36 (†)	26.36(+2.62)
R-2	3	3.99	5.44	4.75 (+0.76)
R-2	All	4.51 (₁)	5.41 (↓)	5.64 (+1.13)
R-L	3	16.57	18.55	17.07 (+0.5)
R-L	All	17.11 (1)	18.65 (↑)	18.83 (+1.72)
BS	3	86.02	86.50	86.25 (+0.23)
BS	All	86.24 (1)	85.96 (1)	86.67 (+0.43)
$BS_{heading}$	3	75.02	85.79	69.50 (-5.52)
BSheading	All	85.37 (1)	85.75 (↓)	85.59 (+0.22)
JS_{text}	3	8.70	10.17	9.12 (+0.42)
JS_{text}	All	9.22 (1)	10.23 (1)	10.38 (+1.16)
CS_{text}	3	32.29	38.35	35.42 (+3.13)
CS_{text}	All	33.23 (↓)	38.74 (38.87 (+5.64)
BLEU	3	0.33	0.5	0.4 (+0.07)
BLEU	All	0.4 (1)	0.5 (1)	0.6 (+0.2)
$CS_{Keyword}$	3	16.20	15.15	15.67 (-0.53)
$CS_{Keyword}$	All	18.14 (↑)	15.42 (1)	16.83 (-2.72)
$JS_{Keyword}$	3	9.85	8.99	9.31 (-0.54)
$JS_{Keyword}$	All	11.10 (†)	9.21 (1)	10.19 (-0.91)
C_{GT}	3	57.65	54.79	45.69 (-12.41)
C_{GT}	All	61.55 (↑)	56.55 (1)	49.43 (-12.12)
C_{LLM}	3	19.76	56.79	76.27 (+57.07)
C_{LLM}	All	40.04 (1)	56.77 (↓)	42.59 (+2.55)
Pr.	3	0.20	0.56	0.65 (+0.46)
Pr.	All	0.4 (1)	0.57 (₁)	0.43 (+0.03)
Rec.	3	0.31	0.55	0.46 (+0.15)
Rec.	All	0.51 (1)	0.56 (₁)	0.46 (-0.05)
F1	3	0.24	0.55	0.54 (+0.3)
F1	All	0.45 (1)	0.56 (₁)	0.44 (-0.01)

Table 12: GPT 40 mini results in "coverage" (Coverage of Ground Truth Limitation (C_{GT}), LLM Generated Limitation (C_{LLM}), Precision, Recall, and F1-score) and "performance" (**R**ouge-x, **B**ert**S**core (over entire text and just heading), Jaccard (JS) and Cosine (CS) similarity over text and heading) metrics when different input sizes and sources are used for limitation extraction. Increasing input length and including external reviews improve generation results (NeurIPS data).

Metric	GPT 40 mini	GPT 40 mini + RAG
Rouge-1	23.69	24.69 (+1.0)
Rouge-2	4.75	5.24 (+ 0.49)
Rouge-L	17.07	18.10 (+1.03)
BERTScore	86.25	86.45 (+0.2)
$BS_{heading}$	69.50	74.74 (+5.24)
JS_{text}	9.12	9.72 (+0.6)
CS_{text}	35.42	34.71 (-0.71)
$CS_{keyword}$	15.67	15.86 (+0.19)
$JS_{keyword}$	9.31	9.60 (+0.29)
BLĖU	0.4	0.7 (+0.3)
C_{GT}	45.69	61.76 (+16.07)
C_{LLM}	76.27	59.18 (-17.09)
Precision	0.65	0.60 (-0.05)
Recall	0.46	0.59 (+0.13)
F1 score	0.54	0.60 (+0.06)

Table 13: Coverage and Performance measure using GPT 40 mini + RAG in NeurIPS data. Ground truth is Limitation + OpenReview, and input is the top 3 sections alongside the abstract.

Prompt 1 = '''
Here is the text containing extracted limitations.
Please identify and list each limitation, ensuring that each one addresses a distinct topic or point.

Figure 4: Prompt to extract limitations from ground truth text.

Metric	Input (3 sections)	Input (All Sections)
Rouge-1	21.32	20.5 (-0.82)
Rouge-2	3.91	4.0 (+0.09)
Rouge-L	14.59	13.95 (-0.64)
BERTScore	84.59	84.23 (-0.36)
JS_{text}	7.70	7.49 (-0.21)
CS_{text}	30.95	31.82 (+0.87)
BLEU	0.3	0.3(0)
$CS_{keyword}$	16.43	17.29 (+0.86)
$JS_{keyword}$	10.00	10.59 (+0.59)
C_{GT}	38.28	37.99 (-0.29)
C_{LLM}	22.09	30.44 (+8.35)
F1 score	0.24	0.32 (+0.08)

Table 14: Ablation study with Mistral 7b in Neurips data. Considering top 3 sections and all sections as input. Author mentioned limitation is ground truth here.

Metric	Input (3 sections)	Input (All Sections)
Rouge-1	33.98	29.57 (-4.41)
Rouge-2	12.02	9.43 (-2.59)
Rouge-L	25.66	22.78 (-2.88)
BERTScore	88.30	87.29 (-1.01)
JS_{text}	14.69	11.16 (-3.53)
CS_{text}	40.40	34.80 (-5.6)
BLEU	3.61	1.22 (-2.39)
$CS_{keyword}$	30.33	30.73 (+0.4)
$JS_{keyword}$	20.28	20.31 (+0.03)
C_{GT}	61.38	18.01 (-43.37)
C_{LLM}	39.04	5.01 (-34.03)
F1 score	0.44	0.05 (-0.39)

Table 15: Ablation study with Llama 3 8b in ACL data. Considering top 3 sections and All sections as input. Author mentioned limitation is ground truth here.

Metric	Limitation	Limitaiton + OpenReview
Rouge-1	21.32	21.02 (-0.3)
Rouge-2	3.91	3.74 (-0.17)
Rouge-L	14.59	14.52 (-0.07)
BERTScore	84.59	84.54 (-0.05)
JS_{text}	7.70	7.62 (-0.08)
CS_{text}	30.95	32.02 (+1.07)
BLEU	0.3	0.2 (-0.1)
$CS_{keyword}$	16.43	14.69 (-1.74)
$JS_{keyword}$	10.00	8.76 (-1.24)
C_{GT}	38.28	32.51 (-5.77)
C_{LLM}	22.09	44.46 (+22.37)
F1 score	0.24	0.34 (+0.1)

Table 16: Ablation study with Mistral 7b in NeurIPS data, considering Limitation and Limitation + OpenReview as ground truth. Top 3 sections are input here.

Approach	Limitation	OpenReview	Limitaiton + OpenReview
Mistral (Traditional) Mistral (Pointwise)	70.67 84.23 (+13.56)	70.19 84.59 (+ 14.4)	70.38 84.54 (+ 14.16)
Llama (Traditional) Llama (Pointwise)	67.55 87.29 (+19.74)	67.17	67.50

Table 17: BERTScore Comparison Between PointWise (Proposed) and Traditional Approaches Across three different Ground Truths (Limitaiton, OpenReview, Limitation + OpenReview) in NeurIPS Data

Agreement	User 1	User 2	User 3
User 1	1	0.75	0.38
User 2	-	1	0.55
User 3	-	-	1

Table 18: Inter annotator agreement amont human annotator for GPT 40 mini as Extractor.

Section	Cosine Similarity
Abstract vs Limitation	33.27
Introduction vs Limitation	33.06
Related Work vs Limitation	25.10
Methodology vs Limitation	26.58
Dataset vs Limitation	25.59
Conclusion vs Limitation	33.04
Experiment and Results vs Limitation	31.73

Table 19: Cosine Similarity between each section and the Limitation section.

Model Name	Top 3 Sections			All Sections		
Widdel Maille	C_{GT}	C _{LLM}	F1 Score	C_{GT}	C_{LLM}	F1 Score
Llama 3	61.38	39.04	0.44	18.01 (-43.37)	5.01 (-34.03)	0.05 (-0.39)
GPT-40 Mini	57.65	19.76	0.24	61.55 (+3.9)	40.28 (+20.52)	0.46 (+0.22)
Mistral	38.28	22.09	0.24	37.99 (-0.29)	30.44 (+8.35)	0.32 (+0.08)

Table 20: Performance (Coverage) of various models considering the top three sections versus the entire paper. Note: C_{GT}, C_{LLM} denotes Coverage of Ground Truth Limitation, Coverage of LLM Generated Limitation. The '+' sign in parentheses indicates the performance gain achieved by using *all sections* instead of only the *top three* sections.

Ground Truth Limitation	LLM Generated Limitation	BS	BS(h)	Match
Error-Prone Models: - Evaluation at the grade school mathematics level indicates that LLMs are still error-prone, but this does not provide insight into their effectiveness in tutoring other subjects.	Reliance on Large Language Models (LLMs): The research paper's dependence on LLMs as Intelligent Tutoring Systems (ITS) for elementary school mathematics may limit the generalizability of the findings. Keywords: tutoring, mathematics, elementary, models, language	87.88	81.08	Yes
Keywords: tutoring, LLMs, mathematics, models, subjects	Errors in Code Generation: The presence of significant errors in code generation indicates that LLMs still face challenges in accurately performing tutoring tasks. Keywords: tutoring, LLMs, errors, code, challenges	90.05	81.38	Yes

Table 21: Ground Truth Limitation matched with multiple LLM Generated Text in PointWise Evaluation. Note: BS, BS(h), Match denotes BERTScore, BERTSCore (heading), and LLM as a judge decision.

Ground Truth Limitation	LLM Generated Limitation	BS	BS(h)	Match
Quality of Representations vs. Range of Words: The primary limitation is not the quality of the word representations produced, but rather the restricted range of words that the methodologies can be applied to. Keywords: representations, words, limitation, methodologies, range	Dictionary Limitations: Dictionaries may not include definitions for all word forms in these languages, which restricts	89.12	83.08	Yes
Dependence on Bilingual Dictionaries : The approach relies on bilingual dictionaries, which may not cover all necessary vocabulary or nuances of the languages. Keywords: bilingual, dictionaries, languages, vocabulary, relies	the applicability of the proposed approach. Keywords: dictionaries, dictionary, languages, forms, definitions	91.71	82.78	Yes

Table 22: LLM Generated Limitation matched with multiple Ground Truth Limitation. Note: BS, BS(h), Match denotes BERTScore, BERTSCore (heading), and LLM as a judge decision.

Ground Truth	Gen. Lim.	GPT 4	Llama 3	User
Model Complexity Concerns: - There is a question regarding whether the performance gains are due to increased model complexity rather than the proposed recursive mixing approach.	**Potential for Increased Complexity**: Although the method is described as simple, the recursive nature of the approach may introduce complexity in implementation and understanding, particularly for practitioners who may not be familiar with the underlying concepts	Yes	Yes	Yes
Computational Complexity: The Dual-aspect Attention mechanism significantly increases computational complexity, which may affect the performance and practicality of the proposed method.	'**Complexity of the Model**: The DACT arch, while innovative, introduces additional complexity compared to traditional models', 'This complexity may lead to longer training times and increased computational resource requirements, which could be a barrier for practical applications in resource-constrained environments'	Yes	Yes	Yes
Insufficient Detailed Experimental Analysis: The paper primarily provides quantitative results without a detailed experimental analysis, which limits the understanding of the findings.	**Limited Discussion on Failure Cases**: While the paper discusses success and failure cases of ex- isting methods, it may not provide sufficient insight into the specific conditions under which SCILL might fail, limiting the understanding of its robust- ness.	Yes	No	Yes
Connection to Practical Methods: There is an unknown connection between the theoretical findings and popular ensemble methods used in practice, which raises questions about the practical relevance of the work.	**Unexplored Variants of Ensemble Methods**: The research does not explore other ensemble methods beyond those mentioned, such as boosting or stacking, which may have different theoretical properties and could provide additional insights into the interpolating regime	No	Yes	Yes
Lack of Exploration of Alternative Methods: The paper does not sufficiently explore or compare the performance of other methods like MFVI and SGLD in the context of covariate shift, which could provide a more comprehensive understanding of the problem.	**Absence of Dataset Details**: The dataset section is marked as "nan," indicating a lack of information about the datasets used for experimentation', 'This omission makes it difficult to assess the validity and applicability of the findings, as the choice of dataset can significantly influence the results	No	No	No

Table 23: LLM as a Judge (GPT 4, Llama 3) and User to decide whether Ground Truth Limitation and LLM Generated Limitation have the same topic or not.

Note: Here GPT 4, Llama 3, and User decision regarding both limitation match or not.

Ground Truth Limitations: Our goal is to provide comprehensive study of learning algorithms, including their hyperparameters. However, our study is limited on the variety of other design choices to restrict the experimental complexity. While we studied both synthetic and realistic image datasets, both these datasets are relatively simple with the same small number of generative factors and each of the factor follows uniform distribution. For learning algorithms, we focus on studying the inductive bias on the representation format while fixing the model arch design which can impact the results. Moreover, we did not study hyperparameters beyond those related to the latent representations. Specifically, we did not study how the type and configurations of the optimizer and the batch size would change the results; instead, we followed common setups in previous studies.

- LLM Refined Limitations: '1. Here are the limitations extracted from the provided limitation:',
- '2. . **Experimental Complexity**: The study is limited in the variety of design choices to reduce experimental complexity.',
- '3.' . **Dataset Complexity**: The study utilizes relatively simple synthetic and realistic image datasets, which have a small number of generative factors.',
- '4. . **Generative Factors**: Each generative factor in the datasets follows a uniform distribution, limiting the diversity of the data.',
- '5. . **Model arch**: The focus is on studying the inductive bias related to the representation format while keeping the model arch design fixed, which may influence the results.', '
- 6. . **Hyperparameter Scope**: The study does not explore hyperparameters beyond those related to latent representations.',
- '7. . **Optimizer and Batch Size**: The impact of different types and configurations of optimizers and batch sizes on the results is not examined, as the study adheres to common setups from previous research.'

Figure 5: Ground Truth Limitations and LLM Extracted Limitations in NeurIPS dataset.

Ground Truth Limitations: While our work provides useful starting point for understanding student feedback, there are limitations to our work. Addressing these limitations will be an important area for future research. Comments may not reflect real student feedback. The comments in our dataset are from users who have chosen to post publicly on youtube. Addi- tionally, the comments may include features specific to this online education setting. Thus, the comments may reflect real student comments from these courses. There is selection bias in lecture sources. Sight includes lectures that may be drawn from the most successful offerings of that course. The instructional quality may not be representative of typical instruction. Thus, inferences drawn about the instruction should be interpreted with caution, as they might not generalize to other lecture settings. We analyze only english comments. We analyze only english comments because the lecture content is given in english and the authors are most comfortable with english. As result, our rubric may not capture the types of feedback from nonenglish students watching lectures taught in english. we annotate small subsample of the data to assess the validity of the automatic labels, we conduct diagnostic study on small, randomly selected subset of the dataset, comprising approximately of the comments. Our work aims to establish preliminary evaluation of the humanmodel agreement and model annotations, and further validation of the automatic labels is necessary.

LLM Refined Limitations: Here are the limitations extracted from the provided limitation:

- 1. **Comments May Not Reflect Real Student Feedback**: The dataset consists of comments from users who have chosen to post publicly on YouTube, which may not accurately represent genuine student feedback.
- 2. **Features Specific to Online Education Setting**: The comments may include characteristics that are unique to the online education context, potentially skewing the interpretation of student feedback.
- 3. **Selection Bias in Lecture Sources**: The lectures included in the analysis may be drawn from the most successful offerings of a course, leading to a selection bias that does not represent typical instructional quality.
- 4. **Generalizability of Instructional Quality**: Inferences about instructional quality based on the selected lectures should be made with caution, as they may not generalize to other lecture settings.
- 5. **Analysis Limited to English Comments**: The analysis is restricted to comments in English, which may overlook feedback from non-English speaking students who are watching lectures taught in English.
- 6. **Small Subsample for Validity Assessment**: The study annotates a small subsample of the data to assess the validity of automatic labels, which may not provide a comprehensive evaluation of the entire dataset
- 7. **Need for Further Validation of Automatic Labels**: The work establishes a preliminary evaluation of human-model agreement and model annotations, indicating that further validation of the automatic labels is necessary.

Figure 6: Ground Truth Limitations and LLM Extracted Limitations in ACL dataset.

Ground Truth OpenReview: Reviewer 1: This paper tries to tackle the domain generalization problem by finding flat minima. It is interesting to investigate the loss landscape of the DNNs on domain generalization problems. Originality: The method proposed in this paper is mainly built on previous approaches such as SWA, and the authors mention this in the submission. Quality: This submission is complete piece of work and the claims are well supported by the empirical results. The evaluation is reasonable. Clarity: This submission is well-written and easy to follow. Significance: This paper provides interesting empirical observations and improves the state-of-the-art by applying their proposed new approach. Pros: The proposed method consistently improves the model performance on domain generalization problems, and the proposed method is flexible and can be easily integrated with other training approaches for better performance. Cons: The generalization bound does not provide intuition for robust risk minimization... Reviewer 2: Originality The paper is, perhaps by design, not particularly original. Rather, it makes the observation that an existing set of techniques with particular goal also has applications to different goal. This is still, to the best of my knowledge, novel observation and contribution, and furthermore the modifications made to SWA to arrive at SWAD also appear novel. Nevertheless, these modifications are relatively small, so the method can still be viewed as version of the general idea of SWA, and thus originality is not the paper' strong suit. Quality The paper is of relatively high quality, primarily due to the empirical results. did not carefully check the theory for correctness, but am unsure as to the significance of the theorems. They appear rather similar to prior theoretical results on how flat minima can lead to better generalization, but with some additional divergence terms thrown in as we are dealing with the DG setting... Reviewer 3: Pros: The paper tackles an important and relevant problem to the ML community. The proposed algorithm is relatively simple to implement in practice and perform well experimentally. Cons: My major concern is wrt to the novelty of the work and the connection/motivation between the proposed algorithm SWAD, domain generalization and the derived theoretical insights...

- LLM Refined OpenReview: "1. . **Connection Between Theory and Method**: The relationship between the theoretical framework presented in the paper and the SWAD algorithm is weak. There is a need for a clearer connection between the flatness term in the equations and the algorithm's performance.",
- '2. . **Novelty of Theoretical Results**: The novelty of the theorem compared to classic results, such as those from Ben-David, is questioned. The theoretical contributions do not significantly advance the understanding of flat minima in the context of domain generalization.',
- '3. . **Generalization Bound Intuition**: The generalization bound provided does not offer sufficient intuition for robust risk minimization (RRM), and it does not convincingly demonstrate that RRM is superior to standard empirical risk minimization (ERM).'],
- '4. . **Dependence on Previous Work**: The proposed method heavily relies on the existing SWA approach, with only minor modifications. This raises concerns about the originality of the contribution.',
- '5. . **Theoretical Framework Limitations**: The reliance on VC dimension in the theoretical framework is problematic, as it can be excessively large for the models used, especially in domain generalization scenarios.',
- '6. . **Overlooked Details in Proofs**: There are concerns that important terms have been overlooked in the technical proofs, which could undermine the validity of the results.',
- '7. . **Lack of Clear Motivation for SWAD**: The motivation for using SWAD over SWA is not clearly articulated, making it difficult to understand the advantages of the proposed algorithm.',
- '8. . **Introduction of New Hyperparameters**: The proposed algorithm introduces additional hyperparameters, which may complicate the implementation and tuning process.',
- '9. . **Limited Societal Impact Discussion**: While a limitations section is included, there is a lack of discussion regarding the societal impacts of the research, although specific concerns were not identified.', "10. . **Potential for Broader Applications**: The method's applicability beyond domain generalization is not explored, which could limit its significance and impact in the broader machine learning community.

Figure 7: Ground Truth and LLM Extracted Limitations in OpenReview.