

Can LLMs Generate Tabular Summaries of Science Papers? Rethinking the Evaluation Protocol

Anonymous ACL submission

Abstract

Literature review tables are essential for summarizing and comparing collections of scientific papers. We explore the task of generating tables that best fulfill a user’s informational needs given a collection of scientific papers. Building on recent work (Newman et al., 2024), we extend prior approaches to address real-world complexities through a combination of LLM-based methods and human annotations. Our contributions focus on three key challenges encountered in real-world use: (i) User prompts are often under-specified; (ii) Retrieved candidate papers frequently contain irrelevant content; and (iii) Task evaluation should move beyond shallow text similarity techniques and instead assess the utility of inferred tables for information-seeking tasks (e.g., comparing papers). To support reproducible evaluation, we introduce ARXIV2TABLE, a more realistic and challenging benchmark for this task, along with a novel approach to improve literature review table generation in real-world scenarios. Our extensive experiments on this benchmark show that both open-weight and proprietary LLMs struggle with the task, highlighting its difficulty and the need for further advancements.

1 Introduction

Literature review tables play a crucial role in scientific research by organizing and summarizing large amounts of information from selected papers into a concise and comparable format (Russell et al., 1993). At the core of these tables are the *schema* and *values* that define their structure, where *schema* refers to the categories or aspects used to summarize different papers and *values* correspond to the specific information extracted from each paper. A well-defined *schema* allows each work to be represented as a row of *values*, enabling structured and transparent comparisons across different studies.

With recent advancements in large language models (LLMs; OpenAI, 2025b; DeepSeek-AI

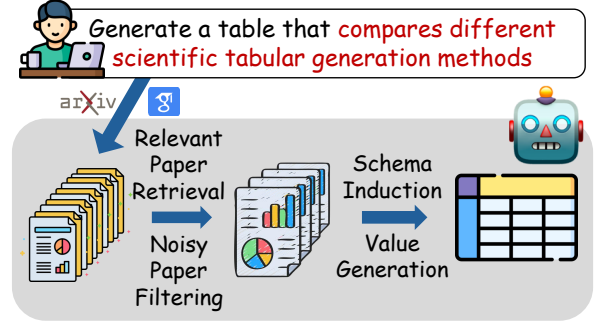


Figure 1: Overview of our proposed task: Given a user’s demand, the LLM first selects the relevant papers that match the request and then generates the schema and values for the desired table.

et al., 2025), several studies (Newman et al., 2024; Dagdelen et al., 2024; Sun et al., 2024) have explored generating literature review tables by prompting LLMs with a set of pre-selected papers and the table’s caption. While these efforts represent meaningful progress, we argue that the existing task definition and evaluation protocols are somewhat unrealistic, thus hindering the practical applicability of generation methods.

First, existing pipelines assume that all provided papers are relevant and should be included in the table. However, in real-world scenarios, distractor papers—those that are irrelevant or contain limited useful information—are common (OpenAI, 2025a). Models should be able to identify and filter out such papers before table construction. Additionally, current pipelines use the ground-truth table’s descriptive caption as the objective for generation. These captions often lack sufficient context, making it difficult for LLMs to infer an appropriate schema, or they may inadvertently reveal the schema and values, leading to biased evaluations.

In this paper, we introduce our task, as illustrated in Figure 1, which improves upon previous task definitions through two key adaptations. First, our pilot study shows that LLMs struggle to retrieve relevant papers from large corpora. To benchmark

this, we introduce distractor papers by selecting them based on semantic similarity to papers in the ground-truth table. LLMs must first determine which papers should be included before generating the table. Second, we replace table captions with abstract user demands that describe the goal of curating the table, making the task more aligned with real-world scenarios. We build upon the ARX-IVDIGESTABLES (Newman et al., 2024) dataset and construct a sibling benchmark through human annotation to verify the selected distractors, comprising 1,957 tables and 7,158 papers.

Meanwhile, current evaluation methods rely on static semantic embeddings to estimate schema overlap between generated and ground-truth tables and require human annotations to assess the quality of unseen schemas and values. However, semantic embeddings struggle to capture nuanced, context-specific variations due to their reliance on pre-trained representations, while human annotation is costly and time-consuming. Moreover, the most effective table generation approaches define schemas primarily based on paper abstracts. This method risks missing important aspects present in the full text, leading to loosely defined schemas with inconsistent granularity.

To address these issues, we propose an annotation-free evaluation framework that instructs an LLM to synthesize QA pairs based on the ground-truth table and assess the generated table by answering these questions. These QA pairs evaluate table content overlap across three dimensions: schema-level, single-cell, and pairwise-cell comparisons. Additionally, we introduce a novel table generation method that batches input papers, iteratively refining paper selection and schema definition by revisiting each paper multiple times. Extensive experiments using five LLMs demonstrate that they struggle with both selecting relevant papers and generating high-quality tables, while our method significantly improves performance on both fronts. Expert validation further confirms the reliability of our QA-synthetic evaluations.

In summary, our contributions are threefold: (1) We introduce an improved task definition for literature review tabular generation, benchmarking it in a more realistic scenario by incorporating distractor papers and replacing table captions with abstract user demands; (2) We propose an annotation-free evaluation framework that leverages LLM-generated QA pairs to assess schema-level, single-cell, and pairwise-cell content overlap, addressing

the limitations of static semantic embeddings and human evaluation; and (3) We develop a novel iterative batch-based table generation method that processes input papers in batches, refining schema definition and paper selection iteratively.

To the best of our knowledge, we are the first to introduce a task that simulates real-world use cases of scientific tabular generation by incorporating user demands and distractor papers, providing a more robust assessment of LLMs in this domain.

2 Related Works

Scientific literature tabular generation Prior works primarily attempt to generate scientific tables through two stages: schema induction and value extraction. For schema induction, early methods like entity-based table generation (Zhang and Balog, 2018) focused on structured input, while recent work has explored schema induction from user queries (Wang et al., 2024) and comparative aspect extraction (Hashimoto et al., 2017). For value extraction, various approaches such as document-grounded question-answering (Kwiatkowski et al., 2019; Dasigi et al., 2021; Lee et al., 2023), aspect-based summarization (Ahuja et al., 2022), and document summarization (DeYoung et al., 2021; Lu et al., 2020) have been proposed to extract relevant information. Beyond these methods, several datasets have been introduced to support scientific table-related tasks, such as TableBank (Li et al., 2020), SciGen (Moosavi et al., 2021), and SciTabQA (Lu et al., 2023). Recently, Newman et al. (2024) proposed streamlining schema and value generation with LLMs sequentially and curated a large-scale benchmark for evaluation. However, all these methods assume a clean and fully relevant set of papers and rely on predefined captions or abstract-based schemas, which risk missing key details. In contrast, we argue for an evaluation approach where candidate papers include tangentially relevant or distracting papers, aligning more closely with real-world literature review workflows.

Table induction for general domains Other than the scientific domain, table induction is also widely studied as text-to-table generation. Prior works attempt this as a sequence-to-sequence task (Li et al., 2023; Wu et al., 2022) or as a question-answering problem (Sundar et al., 2024; Tang et al., 2023). Similar to these works, our framework is capable of better handling both structured and distractive input for real-world literature

review and knowledge synthesis.

3 Task Definition

We first define a pipeline consisting of three sub-tasks that extend prior definitions and better capture the real-world usage of literature review tabular generation. For all the following tasks, we are given a user demand prompt p , which specifies the intended purpose of creating the table. **(T1) Candidate Paper Retrieval:** We begin with a given *universe* of papers (e.g., the content of Google Scholar or arXiv) from which relevant papers need to be identified. Given a large collection, the goal is to use a search engine (IR) to retrieve a subset of *candidate* papers $C := \{d_i\}_{i=1}^M$ of size M , which may include distractor papers—i.e., papers that resemble the user demand prompt but do not fully satisfy the requirement. **(T2) Paper Selection:** Given C , the second subtask is to select the *relevant* subset of size m ($m < M$): $R := \{d_i\}_{i=1}^m \subseteq C$, which best aligns with the user demand p . T2 differs from T1 in scale. Due to the large scale of T1, IR engines must optimize for recall, ensuring that as many relevant papers as possible are retrieved. However, T2 operates at a smaller scale, where precision is the priority, as it focuses on filtering out distractors and selecting only the most relevant papers. **(T3) Table Induction:** Given the selected papers R , the objective is to generate a table with m rows and N columns, where $N \geq 2$ (i.e., no single-column tables). Each row $r_i \in \{r_1, r_2, \dots, r_m\}$ corresponds to a unique input document $d_i \in R$, and each column $c_j \in \{c_1, c_2, \dots, c_N\}$ represents a unique aspect of the documents. We refer to these N columns as the *schema* of the table and the $N \times m$ cells as the *values* of the table. The value of each cell is derived from its respective document according to the aspect defined by the corresponding column.

4 ARXIV2TABLE Construction

We then construct ARXIV2TABLE based on the ARXIVDIGESTABLES dataset which consists of literature tables (extracted from computer science papers) and their corresponding captions. We filter out tables that are structurally incomplete or lack full text for all referenced papers. As a result, we are left with 1,957 tables (with captions) which have rows referring to 7,158 papers. Our construction involves three pillars: user demand inference (§4.1), a simulated paper retrieval (§4.2)

and evaluation through utilization (§4.3).

4.1 Constructing User Demand Prompts

The first step is to collect user demands p that explicitly describe the desired table (can be understood without the table content) and do not reveal the table’s schema or specific values.

Table captions are not appropriate prompts

While the input dataset contains one caption per table, collected from arXiv papers, these captions are meant to complement tables rather than fully describe them. As a result, they are generally concise. For example, a table caption might read: “*Performance comparison of different approaches*,” which is too vague to understand without seeing the table. Consequently, using table captions as prompts may not yield a well-defined task. A more contextually self-contained rewritten user demand might instead be: “*Draft a table that compares different knowledge editing methods, focusing on their performance on QA datasets*.”

Our prompt construction To address this issue, we propose rewriting the captions of literature review tables into abstract yet descriptive user intentions using LLMs. We guide GPT-4o with a prompt (see §A) that first explains the task to the LLM, specifying that the user demand should be sufficiently contextualized to clearly state the table’s purpose while avoiding the inclusion or direct description of column names or specific values. GPT-4o is then expected to infer the user demand for the given table and its caption. For simplicity, we collect only one user demand per table. More examples are provided in Appendix D.

Table captions vs. constructed user demand prompts

To verify that our collected user demands align with our objective, we visualize: (1) the distribution of the number of tokens in the original and modified user demands, and (2) the ratio of captions and user demands of different lengths that have token overlap with the schema or values. From Figure 2, we observe that our modified user demands are generally longer than the original captions, providing a more detailed description of the table’s goal. Furthermore, as shown in Table 1, user demands exhibit a significantly lower overlap ratio with the schema and table values, resulting in fewer overlapping tokens. This ensures a fairer subsequent evaluation.

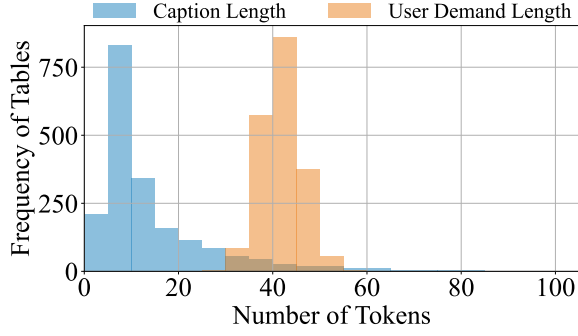


Figure 2: Distribution of the number of tokens between original captions and our modified user demands.

4.2 Paper Retrieval Simulation

The unreliability of paper retrieval Next, we approach the first subtask, candidate paper retrieval, by conducting a pilot study to assess whether LMs can reliably retrieve relevant papers from a large corpus. For each table, we employ a SentenceBERT (Reimers and Gurevych, 2019) encoder as a retrieval engine, selecting papers from the entire corpus based on the highest similarity between the table’s user demand and each paper’s title and abstract. We vary the number of retrieved papers between 2 and 100 and plot the precision and recall of retrieval against the ground-truth papers in the original table (Figure 3).

We observe consistently low precision and recall across different retrieval sizes, highlighting the challenge of retrieving relevant papers from a noisy corpus. This demonstrates that the first subtask is non-trivial and may introduce noise into subtask T2. However, various information retrieval engines, such as Google Scholar and Semantic Scholar, can replace LMs in this subtask. Thus, we decide to simulate T1 by manually adding noisy distractor papers into C to construct R , ensuring a noisy input for T2. This allows us to focus on evaluating LLMs’ capabilities in the T2 and T3 subtasks.

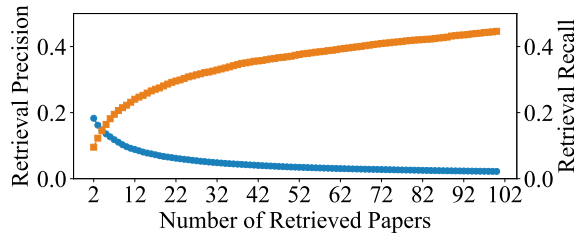


Figure 3: Precision and recall curves for different numbers of retrieved papers.

Similarity-based paper retrieval Moving forward, we associate distractor paper candidates with each table to simulate a potentially noisy document pool before constructing the table. Ideally, distrac-

Prompt	Content	#Table ↓	#Tokens ↓
Caption	Schema	101 (5.2%)	1.2
	Value	46 (2.4%)	1.3
User Demand	Schema	14 (0.7%)	1.0
	Value	8 (0.4%)	1.0

Table 1: Overlap statistics between prompts (the original caption or our constructed user demand) and table content (schema or values). **#Table**: Number (and %) of tables with at least one token from table content overlapping with the prompt. **#Tokens**: Average count of overlapping tokens between table content and prompt.

tor candidates should be semantically related to the table but exhibit key differences that fail to meet the user demand. To select such candidates, we adopt a retrieve-then-annotate approach. First, we use a SentenceBERT encoder F to obtain embeddings for (1) the user demand $F(p)$ and (2) all papers in the corpus $\{F(d_i) \mid d_i \in C\}$. Each paper’s embedding is computed by encoding the concatenation of its title and abstract. We then rank all papers $d_i \notin R$ based on the average of two cosine similarities: (1) the similarity between the candidate and the user demand, and (2) the average similarity between the candidate and each referenced paper:

$$s(d_i) = \cos(F(d_i), F(p)) + \frac{1}{m} \sum_{j=1}^m \cos(F(d_i), F(d_{u_j})).$$

Higher values of $s(d_i)$ indicate stronger semantic relevance, and we select the top 10 ranked papers for each table as its distractor candidates.

Candidates verification via human annotation

After selecting these candidates, we conduct human annotations to verify whether they should indeed be excluded from the table. Given that annotating these tables requires expert knowledge in computer science, we recruit seven postgraduate students with research experience in the field as annotators. To ensure they are well-prepared for the task, the annotators undergo rigorous training, including pilot annotation exams. Their task is to make a binary decision on whether a given distractor paper—based on its title, abstract, user demand, the ground-truth table, and the titles and abstracts of all referenced papers—should be included in the table. Each table contains annotations for 10 papers, with each distractor paper initially assigned to two randomly selected annotators. If both annotators agree on the label, it is finalized. Otherwise, two additional annotators review the paper until a consensus is reached. In the first round, the inter-annotator agreement (IAA) is 94% based on

pairwise agreement, and the Fleiss’ Kappa (Fleiss, 1971) score is 0.73, indicating a substantial level of agreement (Landis and Koch, 1977). Finally, for each table, we randomly select a number of distractor papers between $[m, 10]$ and merge them with R to form C .

4.3 Evaluation via LLM-based Utilization

After constructing the benchmark, we propose evaluating the quality of generated tables from a utilization perspective to address the challenge of aligning schemas and values despite potential differences in phrasing. This is achieved by synthesizing QA pairs based on the ground-truth table and using the generated table to answer them, or vice versa. The flexibility of this QA synthesis allows us to evaluate multiple dimensions of the table while ensuring a structured and scalable assessment. An overview with running examples is shown in Figure 4.

Dimensions of evaluating a table with QAs We introduce three key aspects for evaluating a table in terms of its usability: (1) **Schema**: whether a specific column is included in the generated schema, (2) **Unary Value**: whether a particular cell from the ground-truth table appears in the generated table, and (3) **Pairwise Value**: whether relationships between two cells remain consistent in the generated table.

Recall evaluation We guide GPT-4o in generating these binary QA pairs based on the ground-truth table. For the first two aspects, we generate QA pairs for all columns and cells, whereas for the third aspect, we randomly sample 10 pairs of cells per table and synthesize them into QA pairs. We then prompt GPT-4o to answer these questions based on the generated table, providing yes/no responses. If the answer cannot be found, the model is instructed to respond with “no,” and vice versa for “yes.” The ratio of “yes” answers indicates how well the generated table preserves the schema, individual values, and pairwise relationships. This represents the **recall** of the ground-truth table, measuring how much original information is retained in the generated table.

Precision evaluation To additionally evaluate **precision**, we reverse the process: instead of generating QA pairs from the ground-truth table, we generate them from the generated table and ask another LLM to answer them using the ground-truth table. The precision score reflects how much of the

generated table’s content is actually supported by the original data. By computing the ratio of “yes” answers, we quantify the accuracy of the generated table in reflecting genuine ground-truth information, as well as any additional useful information not present in the ground-truth table.

5 Tabular Generation Methodologies

We explore a range of methods to evaluate on our proposed task, starting from several baselines inspired by prior work (§5.1) and then our proposed approach (§5.2).

5.1 Baseline Methods

We first introduce three methods for generating literature review tables to evaluate their performance on our task and use them as baselines for our proposed method. For easy reference, these methods are termed numerically.

First, **Method 1** generates the table in a one-step process. It takes all available papers R and the user demand p as input, and the model is asked to select all relevant papers and output a table with a well-defined schema and filled values in a single round of conversation. However, this method struggles with extremely long prompts that exceed the LLMs’ context window when generating large tables.

To address this issue, **Method 2** processes papers individually. For each document, the model decides whether it should be included based on the user demand. If included, the model generates a table for that document. After processing all documents, the final table is created by merging the schemas of all individual tables using exact string matching and copying the corresponding values. While this approach reduces the input prompt length, it results in highly sparse tables due to inconsistent schema across papers and the potential omission of relevant information when individual papers lack sufficient context to define comprehensive table aspects.

To overcome both issues, **Method 3** (Newman et al., 2024) introduces a two-stage process. In the first stage, the model selects papers relevant to the user demand based on their titles and abstracts, then generates a corresponding schema. In the second stage, the model loops through the selected papers and fills in the respective rows based on the full text of each document. A minor drawback of this method is that the schema is generated solely from titles and abstracts, which may overlook details present only in the full text.

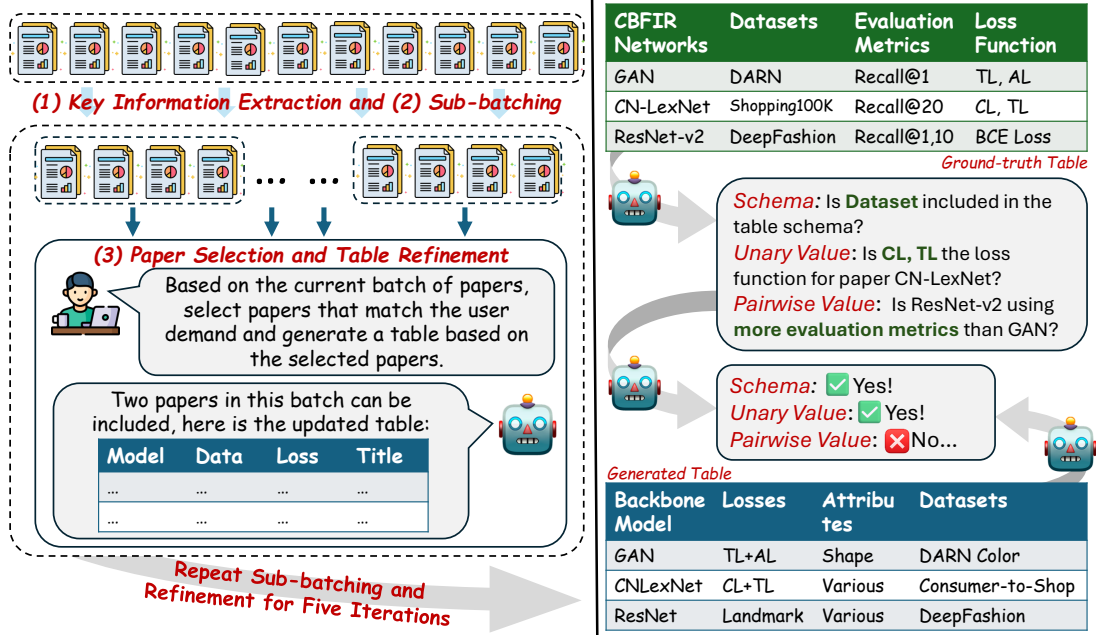


Figure 4: Overview of our proposed iterative batch-based tabular generation method (left) and LLM-based QA-synthesis evaluation protocol (right) with running examples.

5.2 Iterative Batch-based Tabular Generation

Then, we introduce our proposed method for generating literature review tables, as illustrated in Figure 4. Our approach consists of three steps: (A) key information extraction, (B) paper batching, and (C) paper selection and schema refinement, where the latter two steps can be iterated multiple times.

(A) Key Information Extraction Processing multiple papers simultaneously using their full text often results in excessively long prompts that exceed the LLMs’ context window. To address this, we first shorten each paper by instructing the LLM to extract key information from the full text that is relevant to the user’s requirements. Notably, we do not rely solely on the abstract, as important details often appear in the full text but are omitted from the abstract. For each paper, we provide the LLM with its title, abstract, and full text, along with the user’s request, and ask it to generate a concise paragraph that preserves all potentially relevant details. These summary paragraphs serve as condensed representations of the papers for subsequent processing.

(B) Paper Batching Next, we divide all key information paragraphs into smaller batches. Processing too many papers at once negatively affects the model’s performance (as demonstrated by the comparison of Method 1 in Table 2), whereas batching facilitates more efficient comparisons within each batch. For simplicity, we set a batch size of 4 and randomly partition R into $\lceil \frac{|R|}{4} \rceil$ batches.

(C) Paper Selection and Schema Refinement

We initialize an empty schema and table, then sequentially process each batch with the LLM by providing it with the user’s request and summaries of batched papers. The LLM is instructed to (1) decide whether each paper should be included or removed based on its key information and (2) refine the schema based on the current batch of papers. Schema refinement involves adding or removing specific columns or modifying existing values to align with different formats. For new papers that are not deemed suitable for inclusion yet are not in the current table, we also prompt the LLM to insert a new row according to the refined schema. This ensures that the table remains dynamically structured, continuously adapting to new information while maintaining consistency across batches.

Afterward, we iterate steps B and C for k iterations. Here k is a hyper-parameter and we set $k = 5$ in our experiments. The rationale is that multiple iterations allow the schema and table contents to progressively improve, ensuring better alignment with user demands. In each iteration, the batches are newly randomized so that each paper is compared with different subsets, enabling more robust decision-making and reducing bias from specific batch compositions. This iterative refinement also mitigates errors from earlier batches by revisiting and adjusting prior decisions based on newly processed information. After completing all iterations, we individually prompt the LLM to revisit the full

Backbone Model	Method	Paper	Schema			Unary Value			Pairwise Value			Avg
		Recall	P	R	F1	P	R	F1	P	R	F1	
LLAMA-3.3 (70B)	Method 1	52.8	31.3	37.7	34.2	29.6	40.4	34.2	28.4	31.8	30.0	32.8
	Method 2	65.4	26.7	69.3	38.5	17.0	56.8	26.2	11.2	22.5	15.0	26.6
	Method 3	61.9	36.4	40.5	38.3	32.8	44.5	37.8	29.5	30.2	29.8	35.3
	Ours	<u>69.3</u>	<u>41.9</u>	55.4	<u>47.7</u>	<u>43.1</u>	<u>62.6</u>	<u>51.1</u>	<u>36.4</u>	<u>46.9</u>	<u>41.0</u>	<u>46.6</u>
Mistral-Large (123B)	Method 1	54.7	33.1	34.5	33.8	31.6	30.4	31.0	15.5	24.7	19.0	27.9
	Method 2	66.8	27.4	<u>65.0</u>	38.5	22.7	47.4	30.7	17.8	30.7	22.6	30.6
	Method 3	67.9	39.9	41.6	40.7	34.7	46.3	39.7	29.9	35.1	32.3	37.6
	Ours	<u>71.3</u>	<u>45.4</u>	56.7	<u>50.4</u>	<u>43.3</u>	<u>61.5</u>	<u>50.8</u>	<u>42.0</u>	<u>49.2</u>	<u>45.3</u>	<u>48.8</u>
DeepSeek-V3 (685B)	Method 1	57.5	38.7	41.7	40.1	32.5	43.8	37.3	28.7	31.8	30.1	35.8
	Method 2	69.8	34.9	<u>69.0</u>	46.4	27.1	55.5	36.4	25.7	32.7	28.8	37.2
	Method 3	70.9	39.4	44.2	41.7	36.6	49.2	42.0	33.3	36.5	34.8	39.5
	Ours	<u>74.3</u>	<u>39.6</u>	56.9	<u>46.7</u>	<u>47.7</u>	<u>65.2</u>	<u>55.1</u>	<u>40.4</u>	<u>49.8</u>	<u>44.6</u>	<u>48.8</u>
GPT-4o-mini	Method 1	55.9	32.0	35.7	33.7	28.9	39.3	33.3	25.0	31.0	27.7	31.6
	Method 2	68.2	31.5	<u>67.7</u>	43.0	27.7	50.8	35.9	21.6	28.3	24.5	34.5
	Method 3	69.3	40.3	45.9	42.9	38.3	47.5	42.4	35.0	37.8	36.3	40.5
	Ours	<u>72.6</u>	<u>46.5</u>	59.7	<u>52.3</u>	49.0	66.7	56.5	<u>43.5</u>	<u>51.9</u>	<u>47.3</u>	<u>52.0</u>
GPT-4o	Method 1	58.5	35.8	43.2	39.2	36.9	41.8	39.2	29.0	34.7	31.6	36.7
	Method 2	70.2	34.2	68.0	45.5	27.9	56.0	37.2	19.4	33.6	24.6	35.8
	Method 3	71.3	45.0	47.9	46.4	38.7	49.8	43.6	36.9	40.0	38.4	42.8
	Ours	74.6	51.5	59.4	55.2	<u>46.1</u>	66.7	<u>54.5</u>	45.9	55.7	50.3	53.3

Table 2: Tabular evaluation results (%) of five LLMs on the ARXIV2TABLE. The best performances within each backbone are underlined and the best among all backbones are **bold-faced**. Avg refers to averaging three F1 scores.

text of the selected papers to verify the values, thereby completing the tabular generation process.

6 Experiments and Analyses

6.1 Experiment Setup

To demonstrate the generalizability of our method and evaluations, we conduct experiments using two proprietary and three open-source LLMs as backbone model representatives: GPT-4o (OpenAI, 2024b), GPT-4o-mini (OpenAI, 2024a), DeepSeek-V3 (685B; DeepSeek-AI et al., 2024), LLAMA-3.3 (70B; Dubey et al., 2024), and Mistral-Large (123B; Mistral-AI, 2024). We apply all baseline methods and our proposed method to each model and use our evaluation framework to assess the quality of the generated tables based on our benchmark, focusing on four aspects: paper selection (**Paper**), schema content overlap (**Schema**), single-cell value overlap (**Unary Value**), and comparisons across cells (**Pairwise Value**). For paper selection, we use **recall** as the metric to measure the number of ground-truth papers successfully selected. For the latter three tasks, we report precision (P), recall (R), and F1 scores (F1), as explained in §4.3.

6.2 Main Evaluation Results

We report the main evaluation results in Table 2 and summarize our key findings as follows:

(1) **All methods and models struggle to distinguish relevant papers from distractors.** For example, even with their best-performing methods, LLAMA-3.3 and GPT-4o achieve only 65.4% and

71.3% recall on average, respectively. This indicates that a significant number of distractor papers are still being included in the generated tables. Additionally, we observe that processing papers individually or using only abstracts for inclusion decisions yields better performance than concatenating full texts. This suggests that excessively long prompts may weaken LLMs’ ability to make accurate inclusion decisions for each paper.

(2) **Aligning generated schemas with the ground-truth table remains challenging.** Among the baselines, the second method consistently achieves higher recall (e.g., 69.3% with LLAMA-3.3), primarily because it generates a larger number of columns, leading to more overlaps with the ground-truth schema. However, other methods exhibit significantly lower recall, indicating that LLMs still struggle to generate meaningful columns that align well with the ground-truth structure.

(3) **While unary values are well preserved, pairwise comparisons suffer substantial losses.** Most methods, especially our proposed approach, extract unary values with relatively high F1 scores. However, extracting and maintaining pairwise relationships remains challenging. For instance, using LLAMA-3.3, our method achieves a unary F1 score of 51.1 but drops to 41.0 for pairwise values. This trend is consistent across different models, suggesting that while individual entries are correctly identified, capturing the relationships between them remains difficult. The significant gap highlights the challenge of preserving complex relational comparisons within the generated tables.

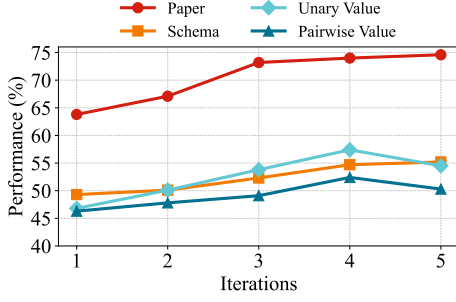


Figure 5: Ablation study on the number of iterations for our iterative batch-based table generation method.

(4) Our proposed method improves performance across all aspects and models. Across all backbone models and evaluation criteria, our method consistently outperforms the baselines. For example, it achieves the highest recall and F1 scores for both unary and pairwise metrics, regardless of model size. This demonstrates that our approach not only enhances overall performance but also provides a more robust solution for handling distractor paper selection and precise table generation.

(5) Larger models lead to better performance. For the three open-source LLMs, we observe a clear trend that increasing the model size improves performance across all aspects when using the same method. For instance, with our approach, scaling from 70B to 123B parameters leads to consistent improvements in most aspects and metrics, reinforcing the importance of stronger generative capabilities in addressing this task.

6.3 Ablation Study on Iteration Number

We further study the impact of the number of iterations, k , in our proposed method to illustrate the importance of refining the schema and table contents over multiple iterations using different batches of papers. As described in §5.2, we perform one round of paper selection and schema refinement five times to achieve optimal performance. In this section, we analyze this process by studying the model’s performance across previous rounds. We select GPT-4o as the backbone model and visualize changes in the recall of paper selection and the F1 scores for schema, unary value, and pairwise comparison overlap by applying the same evaluation protocol to the generated tables after completing iterations ranging from 1 (the first cycle) to 5.

The results are plotted in Figure 5. We observe that during the first four iterations, performance steadily improves across all aspects, demonstrating the effectiveness of iteratively refining paper selection and table schema through multiple itera-

Table	Schema	Unary Value	Pairwise Value
Source	99.5%	100%	98.5%
Target	98.5%	99.5%	97.0%

Table 3: Expert acceptance rate for the synthesized QA pairs sampled from our evaluations.

tions and comparisons between different subsets of papers. At the fifth iteration, however, the improvement slows down, and in some cases, performance even decreases. One possible reason is that the table starts overfitting by including additional values that do not appear in the ground-truth table, reducing precision and leading to lower F1 scores. Considering the overall performance, $k = 5$ is supported as the optimal number of iterations.

6.4 Expert Validation on Synthesized QAs

Lastly, we verify the reliability of synthesizing QA pairs with LLMs for evaluating tabular data. To achieve this, we conduct expert annotations by inviting the authors to manually inspect a random sample of 200 QA pairs covering schema, unary value, and pairwise value comparison aspects. They are asked to annotate (1) whether the generated QA pair is firmly grounded in the source table and (2) whether the LLM correctly answers it based on the target table. The expert acceptance rates are reported in Table 3. We observe that our LLM-synthesized QA pairs are highly reliable, with most acceptance rates above 98% for both source and target tables across schema, unary, and pairwise value comparisons. These results support our evaluation protocol, demonstrating that LLMs can effectively automate the assessment of semantically diverse tabular data.

7 Conclusions

In this work, we introduce an improved literature review table generation task that incorporates distractor papers and replaces table captions with abstract user demands to better align with real-world scenarios, and curated an associated benchmark. Additionally, we propose an annotation-free evaluation framework using LLM-synthesized QA pairs and a novel method to enhance table generation. Our experiments show that current LLMs and existing methods struggle with our task, while our approach significantly improves performance. We envision that our work paves the way for more automated and scalable literature review table generation, ultimately facilitating the efficient synthesis of scientific knowledge in large-scale applications.

Limitations

A minor limitation is that our work uses ARXIVDIGESTABLES as the source of literature review tables for subsequent data reconstruction. However, Newman et al. (2024) have included their pipeline for scalably extracting literature review tables from scientific papers, thus resolving the data reliance gap. Another limitation of our work is its reliance on GPT-4o, a proprietary LLM, for benchmark curation and subsequent evaluation, which may introduce several issues. First, it raises concerns about data contamination (Deng et al., 2024; Dong et al., 2024), as the model may generate user demands (during benchmark curation) and synthesis evaluation questions (when evaluating a generated table against the ground truth) that are similar to its training data, potentially leading to inflated performance in table generation. A data provenance check (Longpre et al., 2024) can be further implemented to address this issue. Second, the benchmark and evaluation process may inherit the internal knowledge or semantic distribution biases of GPT-4o, which could skew the evaluation of other LLMs and reduce the generalizability of our findings. Lastly, a minor issue is scalability, as curating larger datasets using a proprietary model can be resource-intensive and may limit accessibility when extending our framework to other literature or domains. Future work can explore the use of open-source LLMs to replicate the entire process for convenient adaptation to other tabular datasets.

Ethics Statement

The ARXIVDIGESTABLES (Newman et al., 2024) dataset used in our work is shared under the Open Data Commons License, which grants us access to it and allows us to improve and redistribute it for research purposes. Regarding language models, we access all open-source LMs via the Hugging Face Hub (Wolf et al., 2020) and proprietary GPT models through their official API¹. The number of these models, if available, is marked in Table 2. All associated licenses for these models permit user access for research purposes, and we commit to following all terms of use.

When prompting GPT-4o to generate user demands and synthetic QA questions, we explicitly state in the prompt that the LLM should not generate any content that contains personal privacy violations, promotes violence, racial discrimination,

hate speech, sexual, or self-harm contents. We also manually inspect a random sample of 100 data entries generated by GPT-4o for offensive content, and none are detected. Therefore, we believe that our dataset is safe and will not yield any negative or harmful impact.

Our human annotations are conducted by recruiting five graduate-level students who have sufficient experience in data collection for training large language models. They are proficient in English, primarily from Asia, and are paid above the minimum wage in their local jurisdictions. They receive thorough training on the task and are reminded to have a clear understanding of the task instructions before proceeding to annotation. The high level of inter-agreement also confirms the quality of our annotation. The expert annotators have agreed to participate as their contribution to the paper without receiving any compensation.

References

- Ojas Ahuja, Jiacheng Xu, Akshay Gupta, Kevin Horecka, and Greg Durrett. 2022. *ASPECTNEWS: aspect-oriented summarization of news documents*. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2022, Dublin, Ireland, May 22-27, 2022, pages 6494–6506. Association for Computational Linguistics.
- John Dagdelen, Alexander Dunn, Sanghoon Lee, Nicholas Walker, Andrew S Rosen, Gerbrand Ceder, Kristin A Persson, and Anubhav Jain. 2024. Structured information extraction from scientific text with large language models. *Nature Communications*, 15(1):1418.
- Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A. Smith, and Matt Gardner. 2021. *A dataset of information-seeking questions and answers anchored in research papers*. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 4599–4610. Association for Computational Linguistics.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui

¹<https://platform.openai.com/>

745	Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li,	Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu,	808
746	Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang	Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu	809
747	Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L.	Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou,	810
748	Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai	Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun,	811
749	Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai	W. L. Xiao, and Wangding Zeng. 2024. Deepseek-v3	812
750	Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong	technical report . <i>CoRR</i> , abs/2412.19437.	813
751	Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan		
752	Zhang, Minghua Zhang, Minghui Tang, Meng Li,	Chunyu Deng, Yilun Zhao, Xiangru Tang, Mark Ger-	814
753	Miaojun Wang, Mingming Li, Ning Tian, Panpan	stein, and Arman Cohan. 2024. Investigating data	815
754	Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen,	contamination in modern benchmarks for large lan-	816
755	Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan,	guage models . In <i>Proceedings of the 2024 Con-</i>	817
756	Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen,	<i>ference of the North American Chapter of the As-</i>	818
757	Shanghao Lu, Shangyan Zhou, Shanhuang Chen,	<i>sociation for Computational Linguistics: Human</i>	819
758	Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng	<i>Language Technologies (Volume 1: Long Papers),</i>	820
759	Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing	<i>NAACL 2024, Mexico City, Mexico, June 16-21, 2024,</i>	821
760	Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun,	pages 8706–8719. Association for Computational	822
761	T. Wang, Wangding Zeng, Wanxia Zhao, Wen Liu,	Linguistics.	823
762	Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao		
763	Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan	Jay DeYoung, Iz Beltagy, Madeleine van Zuylen, Bailey	824
764	Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin	Kuehl, and Lucy Lu Wang. 2021. Ms^{v2}: Multi-	825
765	Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li,	document summarization of medical studies . In <i>Pro-</i>	826
766	Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin,	<i>ceedings of the 2021 Conference on Empirical Meth-</i>	827
767	Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxi-	<i>ods in Natural Language Processing, EMNLP 2021,</i>	828
768	ang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang,	<i>Virtual Event / Punta Cana, Dominican Republic, 7-</i>	829
769	Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang	<i>11 November, 2021</i> , pages 7494–7513. Association	830
770	Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng	for Computational Linguistics.	831
771	Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi,		
772	Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang,	Yihong Dong, Xue Jiang, Huanyu Liu, Zhi Jin, Bin Gu,	832
773	Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo,	Mengfei Yang, and Ge Li. 2024. Generalization or	833
774	Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yu-	memorization: Data contamination and trustworthy	834
775	jia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You,	evaluation for large language models . In <i>Findings of</i>	835
776	Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu,	<i>the Association for Computational Linguistics, ACL</i>	836
777	Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu,	<i>2024, Bangkok, Thailand and virtual meeting, August</i>	837
778	Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan,	<i>11-16, 2024</i> , pages 12039–12050. Association for	838
779	Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean	Computational Linguistics.	839
780	Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao,	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,	840
781	Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zi-	Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,	841
782	jia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song,	Akhil Mathur, Alan Schelten, Amy Yang, Angela	842
783	Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu	Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang,	843
784	Zhang, and Zhen Zhang. 2025. Deepseek-r1: Incen-	Archi Mitra, Archie Sravankumar, Artem Korenev,	844
785	tivizing reasoning capability in llms via reinforce-	Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien	845
786	ment learning . <i>Preprint</i> , arXiv:2501.12948.	Rodriguez, Austen Gregerson, Ava Spataru, Bap-	846
		tiste Rozière, Bethany Biron, Binh Tang, Bobbie	847
787	DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingx-	Chern, Charlotte Caucheteux, Chaya Nayak, Chloe	848
788	uan Wang, Bochao Wu, Chengda Lu, Chenggang	Bi, Chris Marra, Chris McConnell, Christian Keller,	849
789	Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan,	Christophe Touret, Chunyang Wu, Corinne Wong,	850
790	Damai Dai, Daya Guo, Dejian Yang, Deli Chen,	Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-	851
791	Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai,	lonsius, Daniel Song, Danielle Pintz, Danny Livshits,	852
792	Fuli Luo, Guangbo Hao, Guanting Chen, Guowei	David Esiobu, Dhruv Choudhary, Dhruv Mahajan,	853
793	Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng	Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes,	854
794	Wang, Haowei Zhang, Honghui Ding, Huajian Xin,	Egor Lakomkin, Ehab AlBadawy, Elina Lobanova,	855
795	Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang,	Emily Dinan, Eric Michael Smith, Filip Radenovic,	856
796	Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang,	Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Geor-	857
797	Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie	gia Lewis Anderson, Graeme Nail, Grégoire Mialon,	858
798	Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu,	Guan Pang, Guillem Cucurell, Hailey Nguyen, Han-	859
799	Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean	nah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov,	860
800	Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao,	Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan	861
801	Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang,	Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan	862
802	Mingchuan Zhang, Minghua Zhang, Minghui Tang,	Geffert, Jana Vranes, Jason Park, Jay Mahadeokar,	863
803	Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang,	Jeet Shah, Jelmer van der Linde, Jennifer Billock,	864
804	Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu	Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi,	865
805	Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge,	Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu,	866
806	Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin	Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph	867
807	Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao	Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia,	868

869	Kalyan Vasuden Alwala, Kartikeya Upasani, Kate	all broken: what will it take to fix them? In <i>Forty-</i>	926
870	Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and	<i>first International Conference on Machine Learning,</i>	927
871	et al. 2024. The llama 3 herd of models . <i>CoRR</i> ,	<i>ICML 2024, Vienna, Austria, July 21-27, 2024</i> . Open-	928
872	abs/2407.21783.	Review.net.	929
873	Joseph L Fleiss. 1971. Measuring nominal scale agree-	Xinyuan Lu, Liangming Pan, Qian Liu, Preslav Nakov,	930
874	ment among many raters. <i>Psychological bulletin</i> ,	and Min-Yen Kan. 2023. SCITAB: A challenging	931
875	76(5):378.	benchmark for compositional reasoning and claim	932
876	Hayato Hashimoto, Kazutoshi Shinoda, Hikaru Yokono,	verification on scientific tables . In <i>Proceedings of the</i>	933
877	and Akiko Aizawa. 2017. Automatic generation of	<i>2023 Conference on Empirical Methods in Natural</i>	934
878	review matrices as multi-document summarization	<i>Language Processing, EMNLP 2023, Singapore, De-</i>	935
879	of scientific papers . In <i>Proceedings of the 2nd Joint</i>	<i>cember 6-10, 2023</i> , pages 7787–7813. Association	936
880	<i>Workshop on Bibliometric-enhanced Information Re-</i>	for Computational Linguistics.	937
881	<i>trieval and Natural Language Processing for Dig-</i>	Yao Lu, Yue Dong, and Laurent Charlin. 2020. Multi-	938
882	<i>ital Libraries (BIRNDL 2017) co-located with the</i>	xscience: A large-scale dataset for extreme multi-	939
883	<i>40th International ACM SIGIR Conference on Re-</i>	document summarization of scientific articles . In	940
884	<i>search and Development in Information Retrieval</i>	<i>Proceedings of the 2020 Conference on Empirical</i>	941
885	<i>(SIGIR 2017), Tokyo, Japan, August 11, 2017</i> , vol-	<i>Methods in Natural Language Processing, EMNLP</i>	942
886	ume 1888 of <i>CEUR Workshop Proceedings</i> , pages	<i>2020, Online, November 16-20, 2020</i> , pages 8068–	943
887	69–82. CEUR-WS.org.	8074. Association for Computational Linguistics.	944
888	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-	Mistral-AI. 2024. Large enough . <i>Mistral AI Blog</i> .	945
889	field, Michael Collins, Ankur P. Parikh, Chris Alberti,	Nafise Sadat Moosavi, Andreas Rücklé, Dan Roth,	946
890	Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken-	and Iryna Gurevych. 2021. Scigen: a dataset for	947
891	ton Lee, Kristina Toutanova, Llion Jones, Matthew	reasoning-aware text generation from scientific tables .	948
892	Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob	In <i>Proceedings of the Neural Information Process-</i>	949
893	Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natu-	<i>ing Systems Track on Datasets and Benchmarks 1,</i>	950
894	ral questions: a benchmark for question answering	<i>NeurIPS Datasets and Benchmarks 2021, December</i>	951
895	research . <i>Trans. Assoc. Comput. Linguistics</i> , 7:452–	2021, virtual.	952
896	466.	Benjamin Newman, Yoonjoo Lee, Aakanksha Naik,	953
897	J Richard Landis and Gary G Koch. 1977. The mea-	Pao Siangliulue, Raymond Fok, Juho Kim, Daniel S.	954
898	surement of observer agreement for categorical data.	Weld, Joseph Chee Chang, and Kyle Lo. 2024. Arx-	955
899	<i>biometrics</i> , pages 159–174.	ivdigestables: Synthesizing scientific literature into	956
900	Yoonjoo Lee, Kyungjae Lee, Sunghyun Park, Dasol	tables using language models . In <i>Proceedings of</i>	957
901	Hwang, Jaehyeon Kim, Hong-In Lee, and Moontae	<i>the 2024 Conference on Empirical Methods in Nat-</i>	958
902	Lee. 2023. QASA: advanced question answering	<i>ural Language Processing, EMNLP 2024, Miami,</i>	959
903	on scientific articles . In <i>International Conference</i>	<i>FL, USA, November 12-16, 2024</i> , pages 9612–9631.	960
904	<i>on Machine Learning, ICML 2023, 23-29 July 2023,</i>	Association for Computational Linguistics.	961
905	<i>Honolulu, Hawaii, USA</i> , volume 202 of <i>Proceedings</i>	OpenAI. 2024a. Gpt-4o mini: advancing cost-efficient	962
906	<i>of Machine Learning Research</i> , pages 19036–19052.	intelligence . <i>OpenAI</i> .	963
907	PMLR.	OpenAI. 2024b. Hello gpt-4o . <i>OpenAI</i> .	964
908	Minghao Li, Lei Cui, Shaohan Huang, Furu Wei,	OpenAI. 2025a. Introducing deep research . <i>OpenAI</i>	965
909	Ming Zhou, and Zhoujun Li. 2020. Tablebank: Ta-	Blog .	966
910	ble benchmark for image-based table detection and	OpenAI. 2025b. Openai o3-mini . <i>OpenAI Blog</i> .	967
911	recognition . In <i>Proceedings of The 12th Language</i>	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert:	968
912	<i>Resources and Evaluation Conference, LREC 2020,</i>	Sentence embeddings using siamese bert-networks .	969
913	<i>Marseille, France, May 11-16, 2020</i> , pages 1918–	In <i>Proceedings of the 2019 Conference on Empiri-</i>	970
914	1925. European Language Resources Association.	<i>cal Methods in Natural Language Processing and</i>	971
915	Tong Li, Zhihao Wang, Liangying Shao, Xuling Zheng,	<i>the 9th International Joint Conference on Natural</i>	972
916	Xiaoli Wang, and Jinsong Su. 2023. A sequence-	<i>Language Processing, EMNLP-IJCNLP 2019, Hong</i>	973
917	to-sequence&set model for text-to-table generation .	<i>Kong, China, November 3-7, 2019</i> , pages 3980–3990.	974
918	In <i>Findings of the Association for Computational</i>	Association for Computational Linguistics.	975
919	<i>Linguistics: ACL 2023, Toronto, Canada, July 9-14,</i>	Daniel M. Russell, Mark Stefik, Peter Pirolli, and Stu-	976
920	2023, pages 5358–5370. Association for Computa-	art K. Card. 1993. The cost structure of sensemak-	977
921	tional Linguistics.	ing . In <i>Human-Computer Interaction, INTERACT</i>	978
922	Shayne Longpre, Robert Mahari, Naana Obeng-Marnu,	<i>'93, IFIP TC13 International Conference on Human-</i>	979
923	William Brannon, Tobin South, Katy Ilonka Gero,	<i>Computer Interaction, 24-29 April 1993, Amsterdam,</i>	980
924	Alex Pentland, and Jad Kabbara. 2024. Position:		
925	Data authenticity, consent, & provenance for AI are		

The Netherlands, jointly organised with ACM Conference on Human Aspects in Computing Systems CHI'93, pages 269–276. ACM.

Liangtai Sun, Yang Han, Zihan Zhao, Da Ma, Zhennan Shen, Baocai Chen, Lu Chen, and Kai Yu. 2024. [Sci-eval: A multi-level large language model evaluation benchmark for scientific research](#). In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, pages 19053–19061. AAAI Press.

Anirudh Sundar, Christopher Richardson, and Larry Heck. 2024. [gtbls: Generating tables from text by conditional question answering](#). *CoRR*, abs/2403.14457.

Xiangru Tang, Yiming Zong, Jason Phang, Yilun Zhao, Wangchunshu Zhou, Arman Cohan, and Mark Gestein. 2023. [Struc-bench: Are large language models really good at generating complex structured data?](#) *CoRR*, abs/2309.08963.

Xingbo Wang, Samantha L. Huey, Rui Sheng, Saurabh Mehta, and Fei Wang. 2024. [Scidasynth: Interactive structured knowledge extraction and synthesis from scientific literature with large language model](#). *CoRR*, abs/2404.13765.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, EMNLP 2020 - Demos, Online, November 16-20, 2020*, pages 38–45. Association for Computational Linguistics.

Xueqing Wu, Jiacheng Zhang, and Hang Li. 2022. [Text-to-table: A new way of information extraction](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 2518–2533. Association for Computational Linguistics.

Shuo Zhang and Krisztian Balog. 2018. [On-the-fly table generation](#). In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval, SIGIR 2018, Ann Arbor, MI, USA, July 08-12, 2018*, pages 595–604. ACM.

Appendices

A Implementation Details

In this section, we provide additional implementation details about our benchmark curation and evaluation pipeline, including the prompt we used and the models we accessed.

A.1 Prompts Used

We first introduce the prompt used to construct the ARXIV2TABLE benchmark, as explained in Section 4. The main step involves prompting LLM is to collect user demands that describe the purpose of creating the table while remaining contextually self-contained and not revealing the actual schema or values of the table. We use the following prompt to instruct GPT-4o in generating these user demands.

Given a literature review table, along with its caption, you are tasked with writing a user demand or intention for the creator of this table. The user demand should be written as though you are instructing an AI system to generate the table. Avoid directly mentioning column names in the table itself, but instead, focus on explaining why the table is needed and what information it should contain. You may include a description of the table’s structure, whether it requires detailed or summarized columns. Additionally, infer the user’s intentions from the titles of the papers the table will include. Limit each user demand to 1-2 sentences. Examples of good user demands are: I need a table that outlines how each study conceptualizes the problem, categorizes the task, describes the data analyzed, and summarizes the main findings. The table should have detailed columns for each of these aspects. Generate a detailed table comparing the theoretical background, research methodology, and key results of these papers. You can use several columns to capture these aspects for each paper. I want to create a table that summarizes the datasets used to evaluate different GNN models, focusing on the common features and characteristics found across the papers listed below. The table should have concise columns to

1084	highlight these dataset attributes. Now,	pairs, answer them using the generated	1135
1085	write a user demand for the table below.	table. Provide only "yes" or "no"	1136
1086	The caption of the table is "<CAPTION>".	responses: If the information is present	1137
1087	The table looks like this:	in the generated table, respond with	1138
1088	<TABLE>	"yes." If the information is missing or	1139
1089	The following papers are included in the	different, respond with "no." Your task	1140
1090	table:	is to generate the QA pairs based on the	1141
1091	<PAPER-1> . . . <PAPER-N>	ground-truth table and then answer them	1142
1092	Write the user demand for this table. Do	based on the generated table. Now, begin	1143
1093	not include the column names in the user	by generating the QA pairs.	1144
1094	demand. Write a concise and clear user		
1095	demand covering the function, topic, and	The distribution of number of papers per table	1145
1096	structure of the table with one or two	in ARXIV2TABLE is shown in Figure 6.	1146
1097	sentences. The user demand is:		
1098	Then, for synthesizing QA pairs from a table,	A.2 Evaluation Implementations	1147
1099	we use the following prompt to guide GPT-4o in	We access all open-source LLMs via the Hugging	1148
1100	generating some QA pairs with answers:	Face library (Wolf et al., 2020). The models used	1149
1101		are meta-llama/Llama-3.3-70B-Instruct,	1150
1102	You will evaluate the quality of a	mistralai/Mistral-Large-Instruct-2411,	1151
1103	generated table by comparing it against	and deepseek-ai/DeepSeek-V3.	1152
1104	a ground-truth table. The goal is	For GPT models, we access them via	1153
1105	to assess whether the generated table	the official OpenAI Batch API ² . The mod-	1154
1106	correctly retains the schema, individual	els used are gpt-4o-mini-2024-07-18 and	1155
1107	values, and pairwise relationships. This	gpt-4o-2024-08-06.	1156
1108	is achieved by generating targeted	Note that the DeepSeek model family has a con-	1157
1109	QA pairs based on the ground-truth	text window limit of 64K tokens, whereas the oth-	1158
1110	table and answering them using the	ers have a limit of 128K tokens. The generation	1159
1111	generated table. Step 1: QA Pair	temperature is set to 0.5 for all experiments. All	1160
1112	Generation Based on the Ground-Truth	experiments are repeated twice and the average	1161
1113	Table Generate binary (Yes/No) QA pairs	performance is reported.	1162
1114	focusing on three aspects: Schema	B Method Efficiency Evaluations	1163
1115	QA Pairs: Check whether a specific	In addition to the empirical evaluation of the gen-	1164
1116	column from the ground-truth table	erated tables in Table 2, we also compare the ef-	1165
1117	appears in the generated table schema.	iciency of different methods based on their gen-	1166
1118	Example: Is Dataset included in the	eration success rate and the average number of	1167
1119	table schema? Unary Value QA Pairs:	tokens used per table. The generation success rate	1168
1120	Check whether a specific cell value	refers to the average proportion of tables success-	1169
1121	from the ground-truth table is present	fully generated within the context window limit of	1170
1122	in the generated table. Example: Is	each backbone model. The statistics are reported	1171
1123	CL, TL the loss function for paper	in Table 4. Our observations indicate that while all	1172
1124	CN-LexNet? Pairwise Value QA Pairs:	baseline methods encounter issues with the context	1173
1125	Check whether a relationship between	window limit, our schema induction method effec-	1174
1126	two values remains consistent in the	tively mitigates this problem. Furthermore, our	1175
1127	generated table. Example: Is ResNet-v2	method achieves comparable token usage while	1176
1128	using more evaluation metrics than GAN?	delivering superior performance, highlighting its	1177
1129	For Schema and Unary Value, generate	advantage.	1178
1130	a QA pair for every column and every	C Annotation Details	1179
1131	cell, respectively. For Pairwise Value,	To ensure the high quality of our human annota-	1180
1132	randomly sample 10 pairs per table and	tions, we implement strict quality control measures.	1181
1133	construct the corresponding QA pairs.		
1134	Step 2: Answering QA Pairs Using the		
	Generated Table After generating the QA		

²<https://platform.openai.com/docs/guides/batch>

Method	GSR	#Tokens
Method 1	48.19%	128K
Method 2	98.23%	167K
Method 3	99.71%	110K
Ours	100.0%	118K

Table 4: Comparison of the efficiency of different methods. GSR stands for generation success rate.

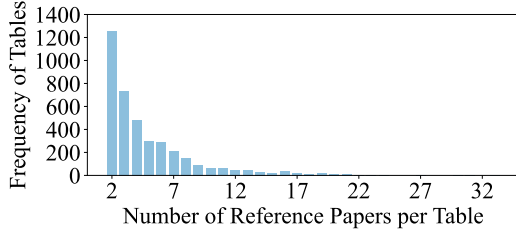


Figure 6: Distribution of number of papers in each table.

First, we select only postgraduate students with research experience in computer science to ensure they are familiar with relevant topics. All selected annotators undergo qualification rounds, and we invite only those who demonstrate satisfactory performance to serve as our main annotators.

For each task, we provide workers with comprehensive task explanations in layman’s terms to enhance their understanding. Additionally, we offer detailed definitions and multiple examples for each choice to help annotators make informed decisions. Each entry requires the worker to provide a binary vote on whether the paper should be excluded or not. Our annotation interface is shown in Figure 7.

To ensure comprehension, we require annotators to confirm that they have thoroughly read the instructions by ticking a checkbox before starting the annotation task. We also manually monitor the performance of annotators throughout the annotation process and provide feedback based on common errors. Spammers or underperforming workers are disqualified. As described in Section 4.2, the inter-annotator agreement supports the quality of our collected annotations.

D Case Studies

Table 5 presents randomly sampled examples of original table captions alongside their improved user demands, demonstrating how refining vague captions enhances specificity and ensures more structured table generation. The findings highlight that well-defined user demands help capture key aspects of table construction, leading to more infor-

mative and targeted tabular representations.

Table 6 illustrates schema, unary value, and pairwise value questions designed to assess the quality of generated tables, ensuring alignment with ground-truth information. The results reveal that this QA-based evaluation effectively quantifies schema retention, individual value accuracy, and consistency in relationships, providing a structured approach for benchmarking table generation models.

Original Table Caption	User Demand
Comparison of Trajectory and Path Planing Approach	Generate a table that compares different trajectory and path planning approaches, focusing on their collision avoidance techniques, benefits, limitations, and applicable scenarios. The table should include detailed columns to capture these aspects for each method mentioned in the relevant papers.
Publications with deep-learning focused sampling methods. We cluster the papers based on the space the sample through and how the samples are evaluated. Some approaches further consider an optional refinement stage.	Create a table that categorizes publications focused on deep-learning-based sampling methods for grasp detection, organizing them by the space in which samples are generated, the evaluation criteria used, and whether a refinement stage is included. The table should provide a comprehensive yet concise overview of the methodological variations and enhancements across different papers.
Categorization of textual explanation methods.	Create a table that categorizes the methods used for providing textual explanations in visual question answering systems, focusing on the types of texts generated and the reasoning processes employed. The table should use succinct columns to differentiate between these methodological aspects for each paper.
Metadata of the three benchmarks that we focus on. XSumSota is a combined benchmark of cite:1400aac and cite:d420ef8 for summaries generated by the state-of-the-art summarization models.	Create a table that details the metadata for three summarization benchmarks, focusing on the composition of annotators, the dataset sizes for validation and testing, and the distribution of positive and negative evaluations. The table should provide a comprehensive comparison across these aspects for each benchmark.
Review of open access ground-based forest datasets	Create a table that reviews various open-access forest datasets, focusing on the publication and data recording years, types of data collected, and their applicability to specific forestry-related tasks. The table should offer a concise summary of each dataset's attributes, including the number of classification categories and geographical location.
Comparison of existing consistency-type models.	Create a table that compares different models focusing on their purpose, the trajectory they follow, the main objects they equate, and their methodological approach. The table should provide detailed insights into how each model addresses consistency issues, drawing from specified papers.

Table 5: Randomly sampled examples of the original captions and their corresponding improved user demands. Most captions are relatively short and may be vague without the full table’s content.

Schema	Unary Value	Pairwise Value
Is Dataset included in the table schema?	Is CL, TL the loss function for paper CN-LexNet?	Is ResNet-v2 using more evaluation metrics than GAN?
Is Model Architecture included in the table schema?	Is GPT-4o the model used for multimodal understanding?	Does GPT-4o have a larger parameter size than LLaMA-2?
Is Training Dataset included in the table schema?	Is ImageNet the dataset used for training ResNet?	Is ResNet trained on more samples than EfficientNet?
Is Performance Metric included in the table schema?	Is BLEU-4 the evaluation metric for MT-BERT?	Does BERT outperform LSTM on BLEU-4 score?
Is Activation Function included in the table schema?	Is ReLU the activation function used in Transformer?	Is GELU smoother than ReLU in function continuity?
Is Optimization Algorithm included in the table schema?	Is Adam the optimizer used for training BERT?	Does Adam converge faster than SGD for BERT training?
Is Pretraining Task included in the table schema?	Is Masked Language Modeling the pre-training task for BERT?	Does BERT use a more complex pretraining strategy than GPT?
Is Hyperparameter included in the table schema?	Is the learning rate set to 0.001 for training ViT?	Does ViT use a higher learning rate than ResNet?
Is Hardware Accelerator included in the table schema?	Is TPU used for training T5?	Do TPUs provide faster training than GPUs for T5?

Table 6: Randomly sampled examples of schema, unary value, and pairwise value questions used to evaluate the quality of generated tables. Each row contains three related questions derived from the same table.

Annotation Task

User Demand

"I need a table that summarizes the key characteristics of various benchmark datasets used in temporal knowledge graph reasoning, including the number of entities, relations, timestamps, and triplets for training, validation, and testing. The table should present this information in a concise manner to facilitate comparison across the studies represented."

Papers in the Current Table

# Entities	# Relations	# Timestamps	# Train Triplets	# Val. Triplets	# Test Triplets
500	20	366	2,735,685	341,961	341,961
15,403	34	198	110,441	13,815	13,800
125,726	203	1,700	323,635	5,000	5,000

Current Literature Review Table

Paper Arxiv Link	Title	Corpus ID
https://arxiv.org/pdf/2104.08419	TIE: A Framework for Embedding-based Incremental Temporal Knowledge Graph Completion	233295959
https://arxiv.org/pdf/1809.03202	Learning Sequence Encoders for Temporal Knowledge Graph Completion	52183483
https://arxiv.org/pdf/2112.05785	TempoQR: Temporal Question Reasoning over Knowledge Graphs	245124416

Paper to Be Decided

Title: Wiki-CS: A Wikipedia-Based Benchmark for Graph Neural Networks

Abstract: We present Wiki-CS, a novel dataset derived from Wikipedia for benchmarking Graph Neural Networks. The dataset consists of nodes corresponding to Computer Science articles, with edges based on hyperlinks and 10 classes representing different branches of the field. We use the dataset to evaluate semi-supervised node classification and single-relation link prediction models. Our experiments show that these methods perform well on a new domain, with structural properties different from earlier benchmarks. The dataset is publicly available, along with the implementation of the data pipeline and the benchmark experiments, at this [https URL](https://arxiv.org/pdf/2007.02901).

Link: <https://arxiv.org/pdf/2007.02901>

Decision

Based on the user demand and the existing literature review table, should this paper be included?

☐ Include
 ☐ Exclude

Submit

Figure 7: The annotation interface we used for collecting the gold labels for distractor papers.