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

Weiqi Wang, Jiefu Ou,

Yangqiu Song, Benjamin Van Durme, Daniel Khashabi

*Under Review at ACL 2025*

# Task: Generating Tabular Summary for Scientific Content

- **Input:** A user prompt seeking scientific information.

- **Output:** A table that summarizes the information/insights extracted from the relevant science papers.

**Prompt:** Generate a table comparing video dataset papers



Repositories of science: arXiv, Google Scholar, etc.

Potentially relevant papers

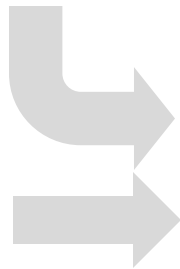| | Dataset size | Annotation method | Intended Application | Evaluation Metric |
|---|---|---|---|---|
| Paper 1 | 1,200 video sequences | Subjectively annotated | Objective VQA method development | Subjective Mean Opinion Score |
| Paper 2 | 585 videos | Subjective video quality scores via crowdsourcing | NR video quality prediction advancement | Subjective video quality scores |
| Paper 3 | 153,841 videos | Coarsely annotated set with five quality ratings each | Deep-learning VQA model training | Spearman rank-order correlation coefficient |
| Paper 4 | 1 million YouTube videos | N/A | Large-scale video classification and action recognition | Performance improvements over baselines |

Tabular summary

# Why solve this task?

- Tables are widely accepted format for scientific content.
    - Improves clarity, enables comparisons
- Integrating it into our workflow will improve quality and interpretability.

**Prompt:** Generate a table comparing video dataset papers



Repositories of science: arXiv, Google Scholar, etc.

Potentially relevant papers

| | Dataset size | Annotation method | Intended Application | Evaluation Metric |
|---|---|---|---|---|
| Paper 1 | 1,200 video sequences | Subjectively annotated | Objective VQA method development | Subjective Mean Opinion Score |
| Paper 2 | 585 videos | Subjective video quality scores via crowdsourcing | NR video quality prediction advancement | Subjective video quality scores |
| Paper 3 | 153,841 videos | Coarsely annotated set with five quality ratings each | Deep-learning VQA model training | Spearman rank-order correlation coefficient |
| Paper 4 | 1 million YouTube videos | N/A | Large-scale video classification and action recognition | Performance improvements over baselines |

Tabular summary

# How good is the state-of-the-art at this task?

ChatGPT

Create a table where the rows are language models:

GPT-3
BLOOM
GPT-4
Gemini
DeepSeek (V3)
DeepSeek (R1)
LLaMA3
Qwen2
OLMo 2

The columns are different details about their alignment (after pre-training) stage.
For example, what alignment algorithm they used, in what order, where did they get their alignment data, etc.

Compiling detailed information about the alignment processes for various language models can be challenging due to the proprietary nature of some models and the evolving practices in the field. Below is a table summarizing the available information on the alignment stages of the specified models:

| Language Model | Alignment Algorithm(s) Used | Alignment Data Sources | Alignment Process Details |
|---|---|---|---|
| GPT-3 | Not publicly disclosed | Not publicly disclosed | Limited information available on post-training alignment methods. |
| BLOOM | Not publicly disclosed | Not publicly disclosed | Limited information available on post-training alignment methods. |
| GPT-4 | Not publicly disclosed | Not publicly disclosed | Limited information available on post-training alignment methods. |
| Gemini | Not publicly disclosed | Not publicly disclosed | Limited information available on post-training alignment methods. |
| DeepSeek (V3) | Not publicly disclosed | Not publicly disclosed | Limited information available on post-training alignment methods. |
| DeepSeek (R1) | Not publicly disclosed | Not publicly disclosed | Limited information available on post-training alignment methods. |

Example retrieved on 02/27/2025

# Recent related work: ArxivDigesTables

- Extracted 2.2K tables from existing papers.

- The table captions serve as the task prompts.

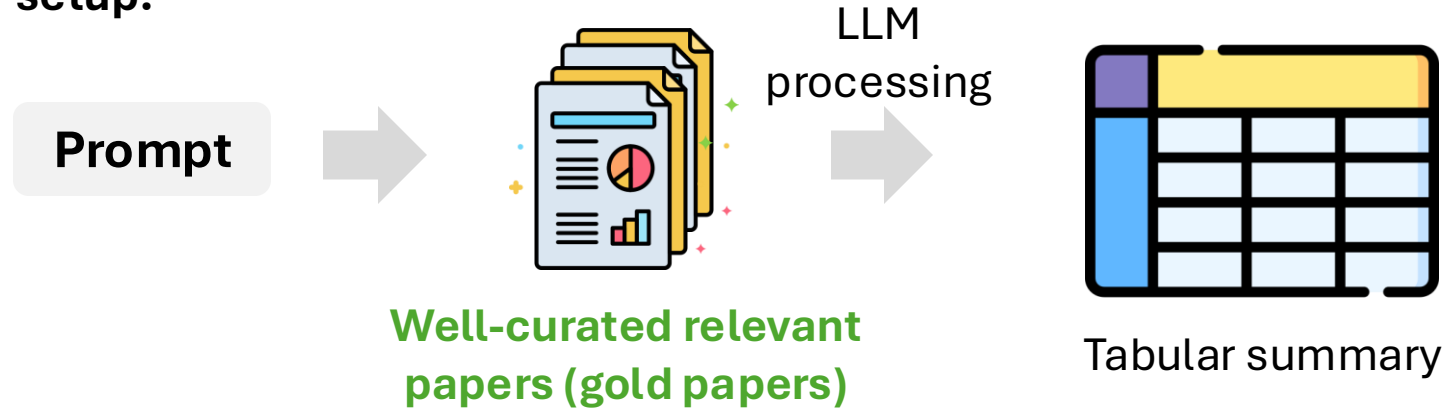- Rows of the table correspond to individual papers (7K) papers.

| | Dataset | Size | Task | Annotations |
|---|---|---|---|---|
| **Paper 1** | KoNViD-1k | 1200 | VQA | 114 |
| **Paper 2** | LIVE-VQC | 585 | VQA | 240 |
| **Paper 3** | KoNViD-150k | 153,841 | VQA | 5 |
| **Paper 4** | Sports-1M | 1,133,158 | Classification | - (auto) |

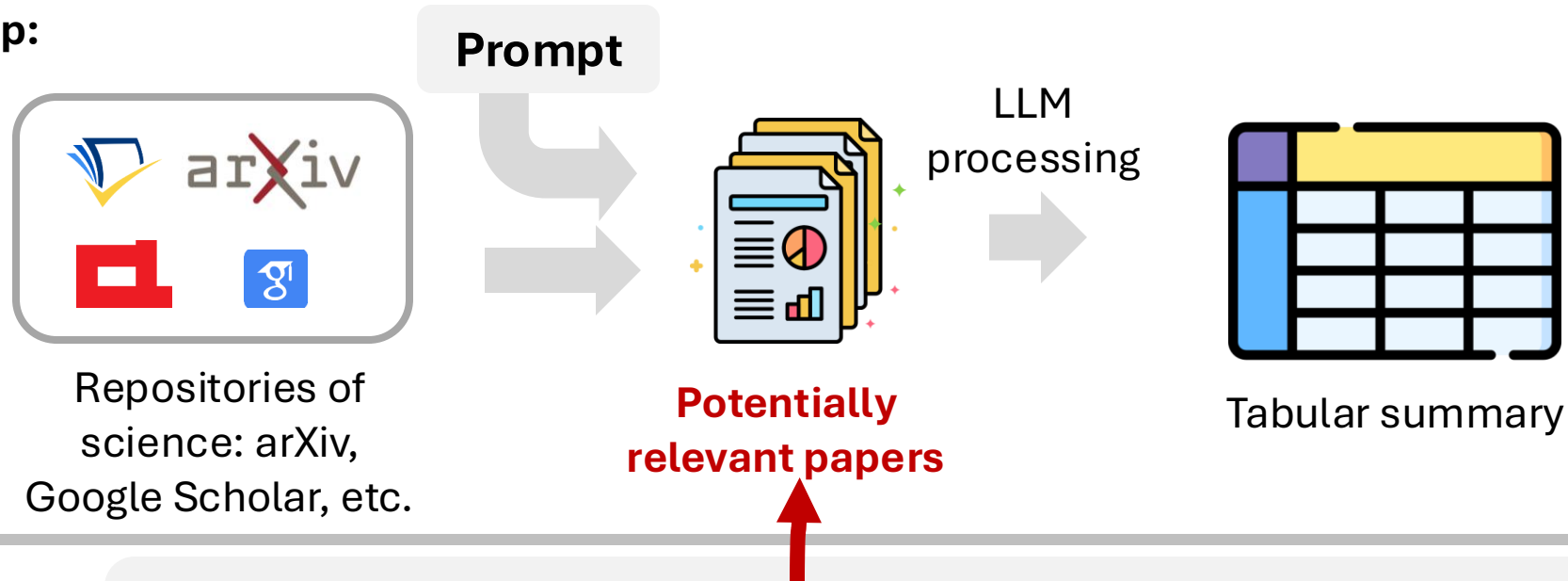**We build upon this work by addressing their weaknesses!**

# Limitations of prior work

1.  **The assumption that papers a carefully curated relevant papers are available, is idealistic in realistic scenarios.**

2.  Table captions are not appropriate task prompts.

3.  Rely on static embedding and human annotation to evaluate generated tables.

**Newman et al. setup:**

Prompt → Well-curated relevant papers (gold papers) → LLM processing → Tabular summary

**Our setup:**

Prompt

Repositories of science: arXiv, Google Scholar, etc. → Potentially relevant papers → LLM processing → Tabular summary

We build a retrieval engine over papers and identify hard negative candidate papers to make evaluation realistic.

# Limitations of prior work

1. The assumption that papers a carefully curated relevant papers are available, is idealistic in realistic scenarios.

2. **Table captions are not appropriate as task prompts.**

3. Rely on static embedding and human annotation to evaluate generated tables.

# User Demand vs. Captions

- Prompts in prior work [Newman et al.] are table captions.

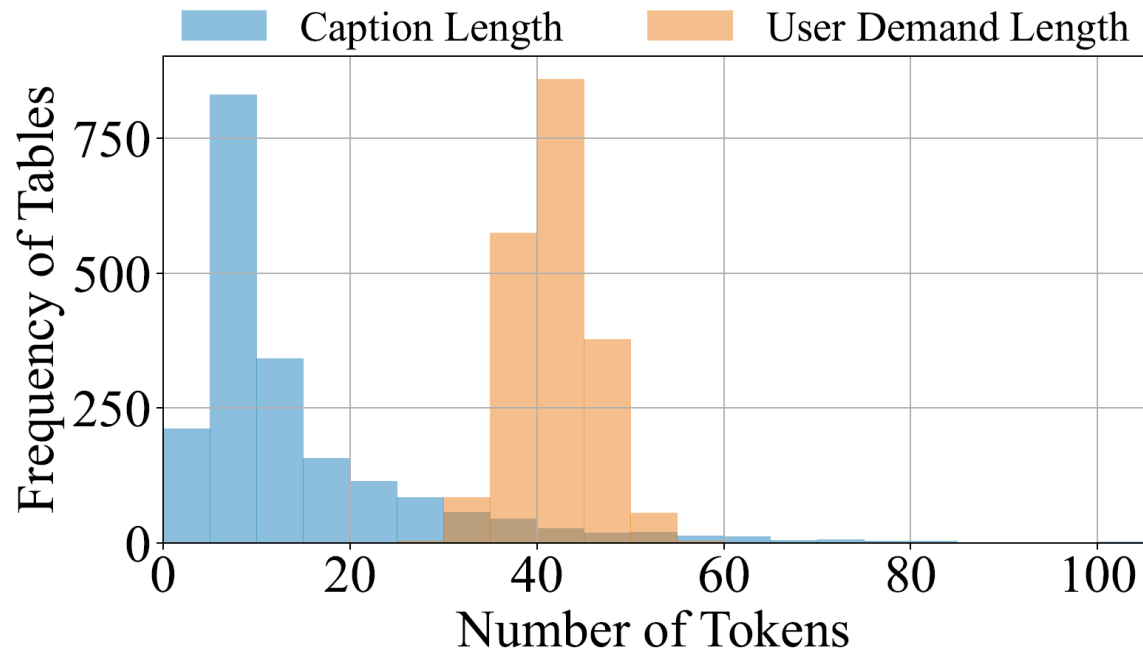  *Comparison of Trajectory and Path Planning Approach*

  Brief and ambiguous

- We replace them with **user demand** prompts:

  *Generate a table that compares different trajectory and path planning approaches. You can focus on their collision avoidance techniques, benefits, limitations, and applicable scenarios.*

  Longer and more precise

  - We collect these by careful prompting of LLMs to:
    - (1) obtain a more complete prompt while
    - (2) avoiding leakage of table schema/values.

# User Demand vs. Captions



Our collected user demands feature longer context, thus including better hints to curate the table.

# Limitations of prior work

1. The assumption that papers a carefully curated relevant papers are available, is idealistic in realistic scenarios.

2. Table captions are not appropriate as task prompts.

3. **Rely on static embedding and human annotation to evaluate generated tables.**

Originally, we have a ground-truth table extracted from a paper

| CBFIR Networks | Datasets | Evaluation Metrics | Loss Function |
|---|---|---|---|
| GAN | DARN | Recall@1 | TL, AL |
| CN-LexNet | Shopping100K | Recall@20 | CL, TL |
| ResNet-v2 | DeepFashion | Recall@1,10 | BCE Loss |

*Ground-truth Table*

We synthesize QA pairs from the ground-truth table about schema and values

table **schema**:
*Is **Dataset** included in the table schema?*

**unary** *(cell) values*:
*Is **CL, TL** the loss function for paper CN-LexNet?*

**pairwise** *comparisons*:
*Is ResNet using **more evaluation metrics** than GAN?*

Then, we ask an LLM to answer these QAs based on the generated table

The ratio of "Correct" indicates the **recall**.

Correct! ✅          Correct! ✅          Incorrect! ❌

Then, we have a table that is generated by an LLM

| Backbone Model | Losses | Attributes | Datasets |
|---|---|---|---|
| GAN | TL+AL | Shape | DARN Color |
| CNLexNet | CL+TL | Various | Consumer-to-Shop |
| ResNet | Landmark | Various | DeepFashion |

*Generated Table*

12

Similarly, we can reverse the process by starting with the generated table.

| Backbone Model | Losses | Attributes | Datasets |
|---|---|---|---|
| GAN | TL+AL | Shape | DARN Color |
| CNLexNet | CL+TL | Various | Consumer-to-Shop |
| ResNet | Landmark | Various | DeepFashion |

*Generated Table*

Again, we synthesize QAs based on the generated table.

table **schema**:
*Is Attributes included in the table schema?*

*unary (cell) values*:
*Is DARN Color used in GAN?*

*pairwise comparisons*:
*Is ResNet using fewer losses than GAN?*

The ratio of "Correct" indicates the **precision**.

But answer them using the ground-truth table.

Incorrect! ❌

Incorrect! ❌

Correct! ✅

And answer QAs using our ground-truth table.

| CBFIR Networks | Datasets | Evaluation Metrics | Loss Function |
|---|---|---|---|
| GAN | DARN | Recall@1 | TL, AL |
| CN-LexNet | Shopping100K | Recall@20 | CL, TL |
| ResNet-v2 | DeepFashion | Recall@1,10 | BCE Loss |

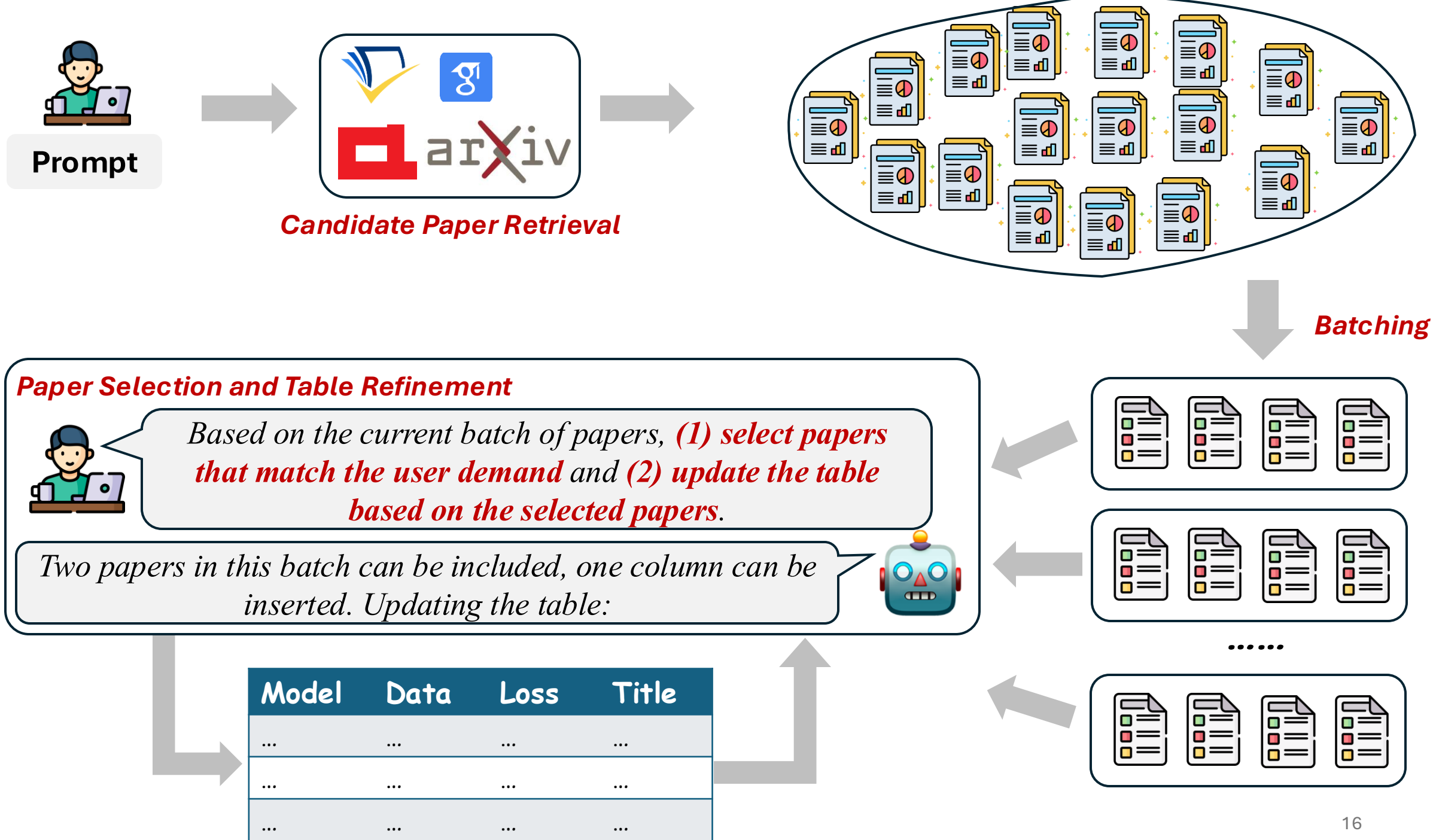*Ground-truth Table*

13

# Our released data: **arXiv2Table**

- Expanded version of Newman et al. 2024.

- Contains
  - 2.1K user demand prompts
  - 2.1K tables (inherited from arXivDigestable).
    - Dropped few low-quality tables.
  - Each prompt comes with it a set of candidate (distractor + gold) papers.
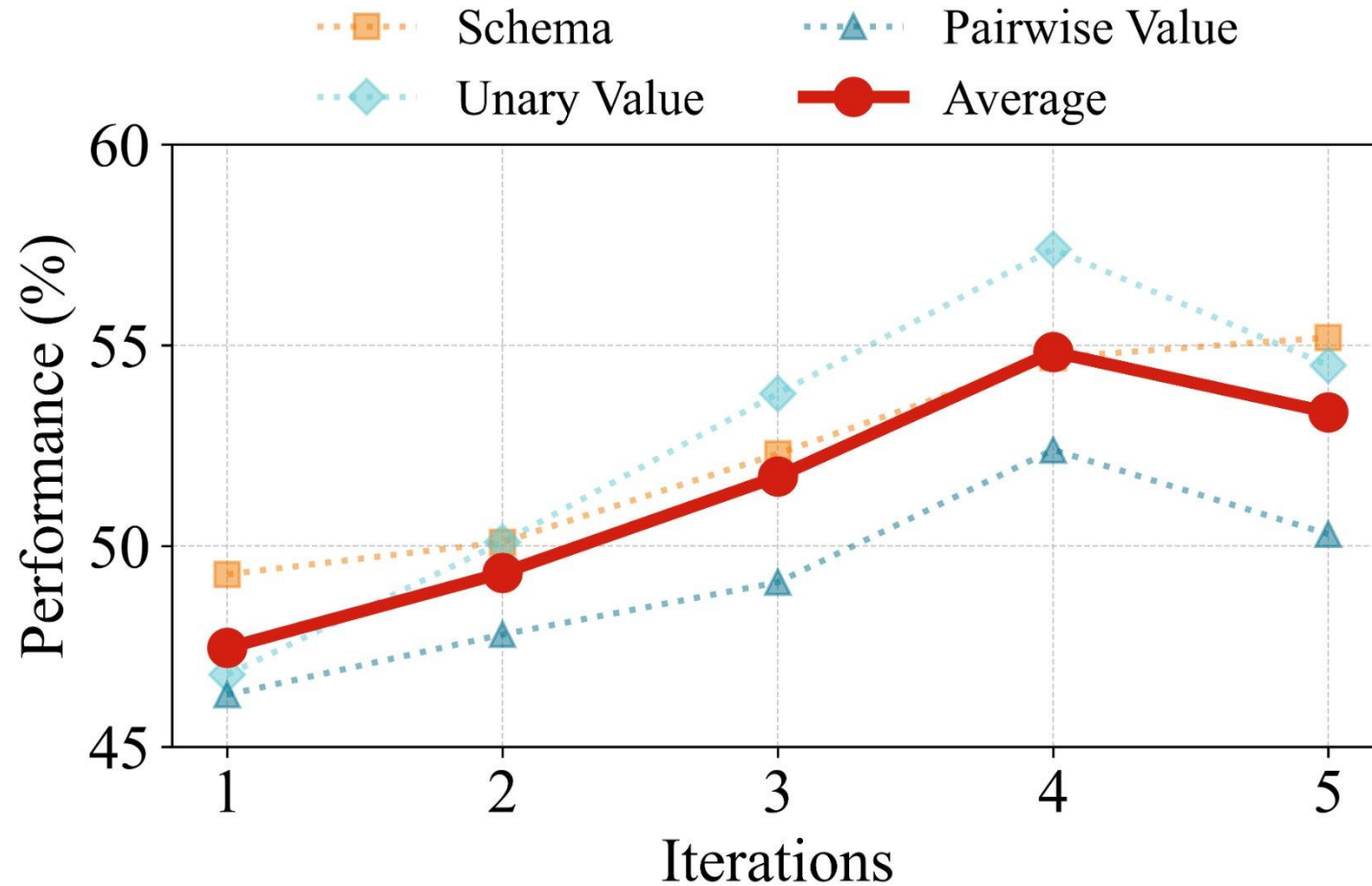  - Evaluation framework based on utilization.

Dataset will be on arXiv on coming weeks!

# We also proposed a new approach

- An inference-time algorithm that iteratively digests and organizes papers in tabular form.
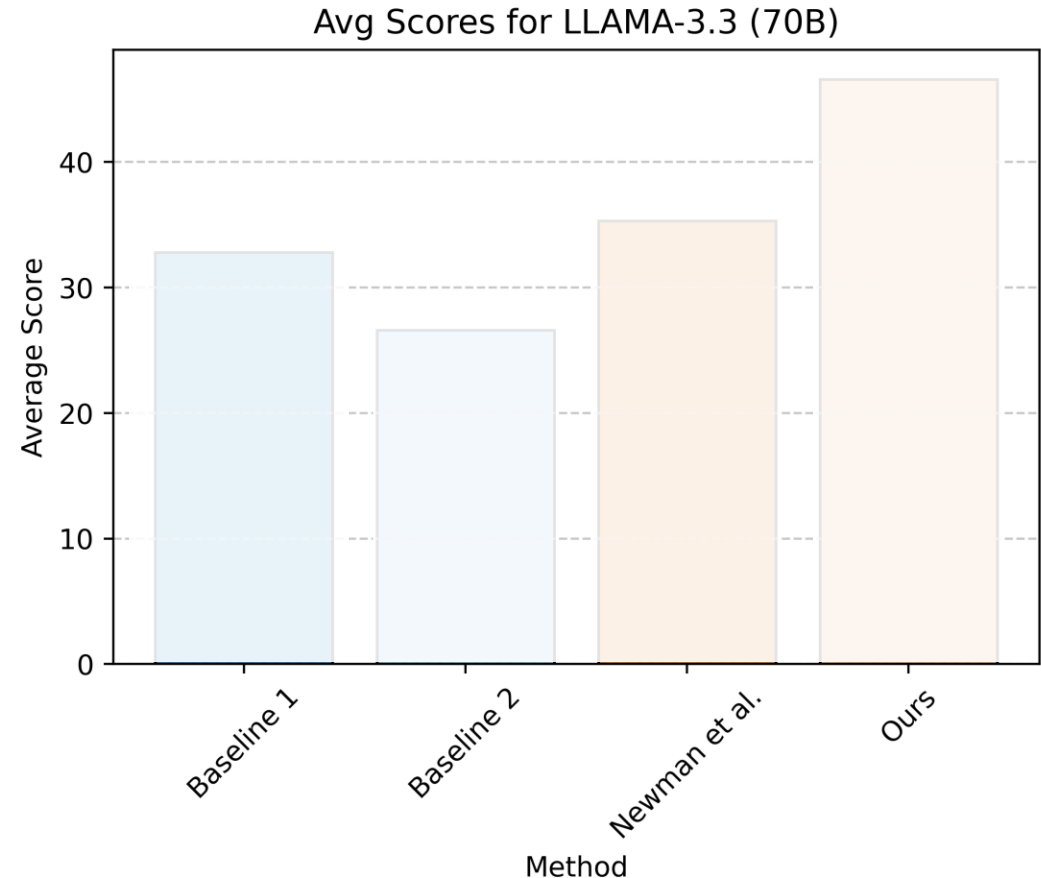
**Prompt**

**Candidate Paper Retrieval**

**Batching**

**Paper Selection and Table Refinement**

Based on the current batch of papers, *(1) select papers that match the user demand* and *(2) update the table based on the selected papers*.

Two papers in this batch can be included, one column can be inserted. Updating the table:

| Model | Data | Loss | Title |
|-------|------|------|-------|
| ... | ... | ... | ... |
| ... | ... | ... | ... |
| ... | ... | ... | ... |

......

# Evaluation vs Number of Iterations



With more iterations, all aspects of the generated tables improves (up to iter ~4).

# Evaluation of the end-to-end pipeline

- **Model:** Llama 3.3 (70B)

- **Baseline 1:** All papers processed in one conversation round.

- **Baseline 2:** One-conversation round per paper.

- **Newman et al.:** Two stages; define schema in the first round, then fill in the values.
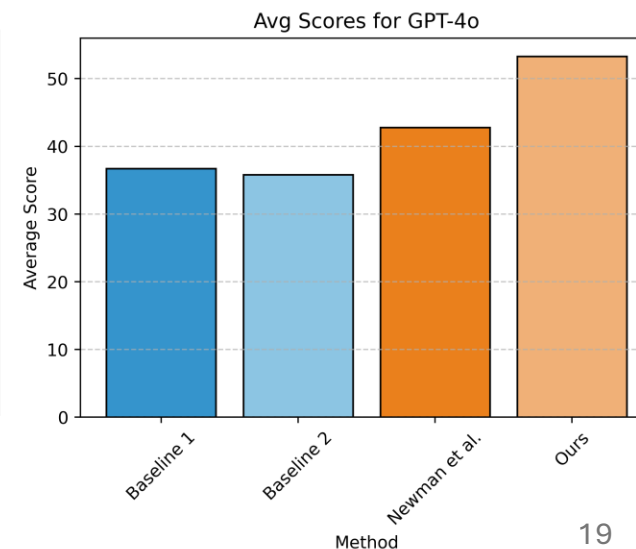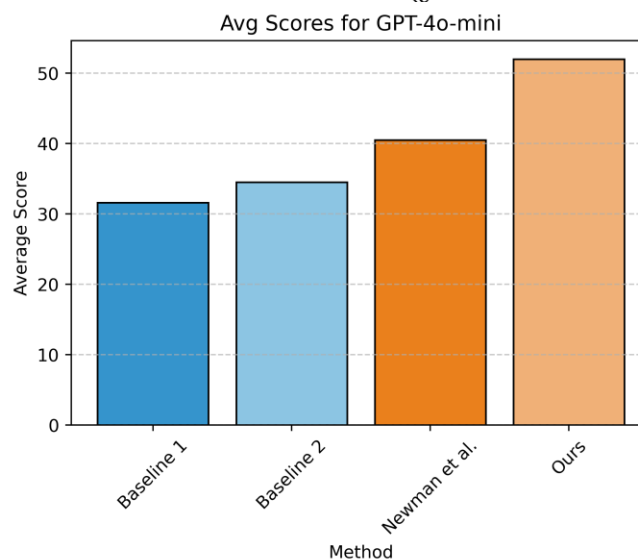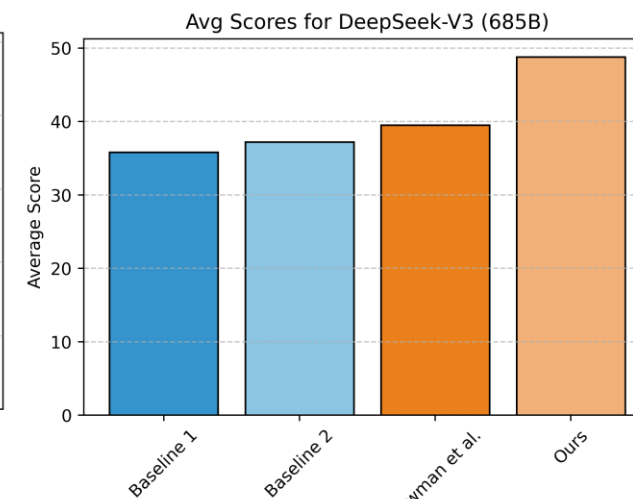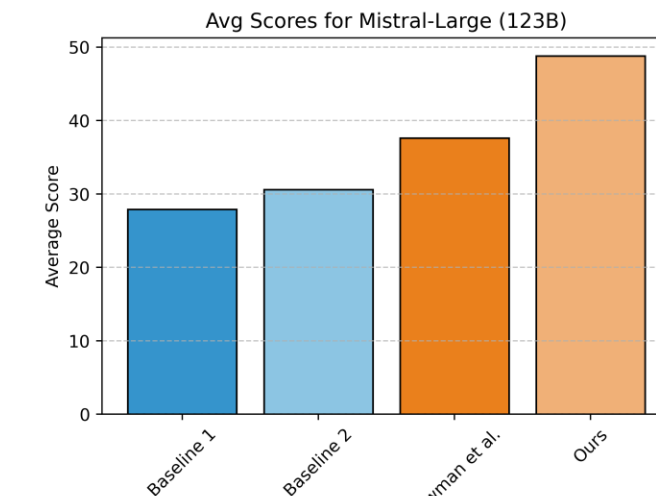


Avg Scores for LLAMA-3.3 (70B)

Our proposed approach outperforms existing results.

# Evaluation of the end-to-end pipeline (other models)

The gains of our approach is consistent across different models.

The task remains challenging for all these approaches.

# Summary and Conclusion

- **Motivation:** A more realistic pipeline for evaluating tabular summarization of science literature.
  - **Why?** Tabular summaries are crucial framework for quickly aggregating and understanding the progress in science.

- We introduce arXiv2Table, a framework for evaluating systems for tabular summarization.

- We also develop a system to address the challenge posed.

- Our benchmark is challenging! Give it a try!! 🚀