

Topic 4. Text-based Content and Comparative Table Generation



Yufang Hou

Professor
Interdisciplinary
Transformation University
yufang.hou@it-u.at

Agenda

- ❑ **Short Text Generation with Citation Grounding**
- ❑ **Automatic Research Paper Writing**
- ❑ **Automated Survey & Deep Research Generation**
- ❑ **Meta-analysis Table Generation and Comparative Literature Synthesis**
- ❑ **Summary**

Cite-worthiness Detection / Citation Recommendation

- **Cite-worthiness detection**: given a sentence in a draft manuscript, classify whether it needs a citation [1].

Should there be a [?] here?

_____ [?]

Cite-worthy?

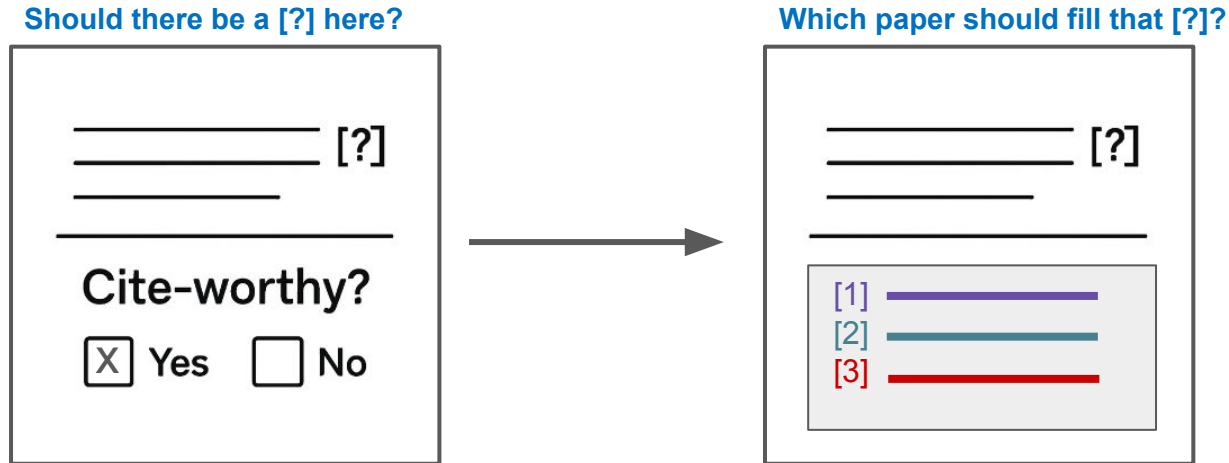
Yes No

[1] Wright & Augenstein, ACL 2021 Findings. [CiteWorth: Cite-Worthiness Detection for Improved Scientific Document Understanding.](#)

[2] Çelik & Tekir, EMNLP 2025. [CiteBART: Learning to Generate Citations for Local Citation Recommendation.](#)

Cite-worthiness Detection / Citation Recommendation

- **Cite-worthiness detection**: given a sentence in a draft manuscript, classify whether it needs a citation [1].
- **Citation recommendation**: given a context that needs a citation (a sentence, paragraph, or full manuscript), retrieve and rank candidate papers to cite, or directly generate citations [2].

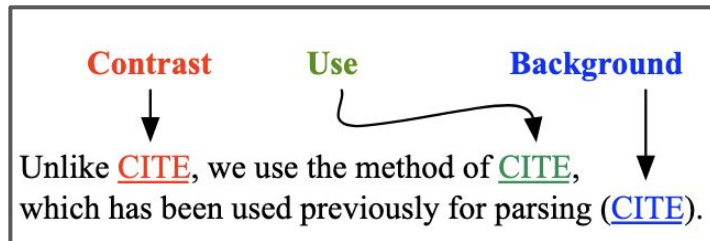


[1] Wright & Augenstein, ACL 2021 Findings. [CiteWorth: Cite-Worthiness Detection for Improved Scientific Document Understanding.](#)

[2] Çelik & Tekir, EMNLP 2025. [CiteBART: Learning to Generate Citations for Local Citation Recommendation.](#)

Citation Analysis: Citation Function

- **Citation function**: classify how scholars use and frame citations [1].



Class	Description	Example
BACKGROUND	<i>P</i> provides relevant information for this domain.	This is often referred to as incorporating deterministic closure (Dörre, 1993).
MOTIVATION	<i>P</i> illustrates need for data, goals, methods, etc.	As shown in Meurers (1994), this is a well-motivated convention [...]
USES	Uses data, methods, etc., from <i>P</i> .	The head words can be automatically extracted [...] in the manner described by Magerman (1994).
EXTENSION	Extends <i>P</i> 's data, methods, etc.	[...] we improve a two-dimensional multimodal version of LDA (Andrews et al, 2009) [...]
COMPARISON OR CONTRAST	Expresses similarity/differences to <i>P</i> .	Other approaches use less deep linguistic resources (e.g., POS-tags Szymne (2008)) [...]
FUTURE	<i>P</i> is a potential avenue for future work.	[...] but we plan to do so in the near future using the algorithm of Littlestone and Warmuth (1992).

Citation Analysis: CORWA

- **CORWA**: citation oriented related work annotation, which decomposes the related work section with three inter-related annotation tasks [1].

1 [Transition] [BOS] Automatic Related work generation is a challenging task.

2 [Narrative_cite] [BOS] Early studies take the extractive approach (Hoang and Kan, 2010; Hu and Wan, 2014).

3 [Transition] [BOS] Recent works switch their attention to the abstractive approach.

4 [Single_summ] [BOS] Xing et al. (2020)

5 extends pointer-generator network (See et al., 2017) to recover a citation sentence given its neighbor sentences and the cited paper's abstract.

6 [Multi_summ] [BOS] While Chen et al. (2021)

7 proposes a custom relation-aware multi-document encoder; Ge et al. (2021)

8 develops a model with multiple inputs and multiple training objectives.

9 [Reflection] [BOS] Although modeling is essential for related work generation, we focus on developing a dataset for related work generation in this work.

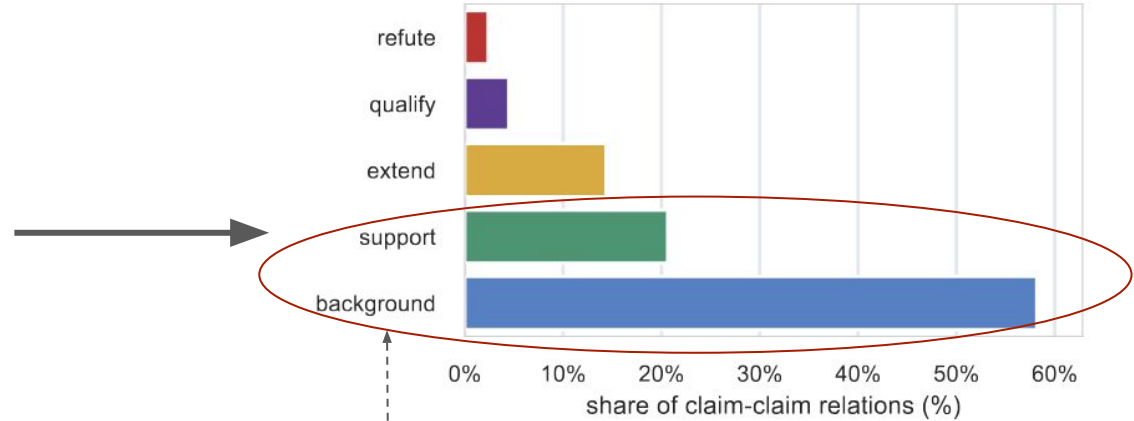
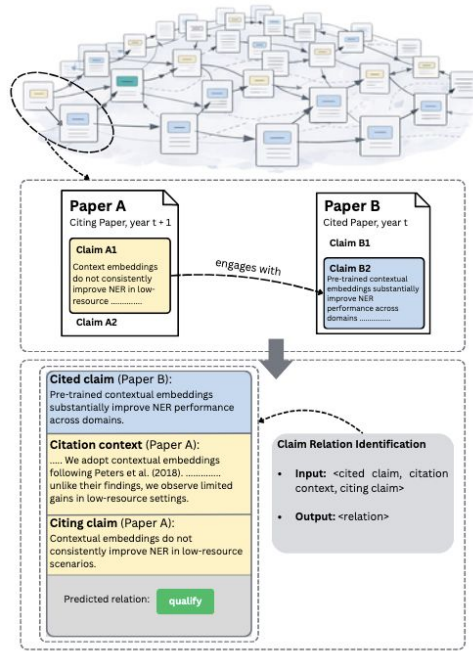
Discourse tagging task tags the role of each related work sentence with one of six labels: {single_summ, multi_summ, narrative_cite, reflection, transition, other}.

Citation span detection task identifies the span of text whose information is directly derived from a specific cited paper.

Citation type recognition task indicates whether a cited work is discussed in detail or used to illustrate a high-level concept.

Citation Analysis: ClaimFlow

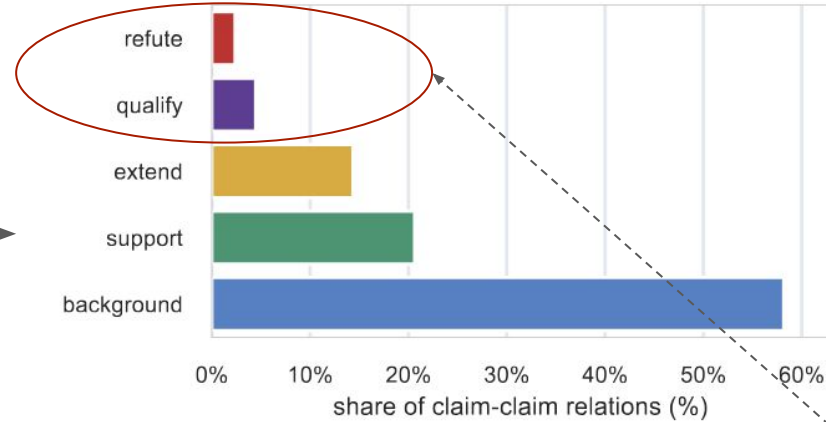
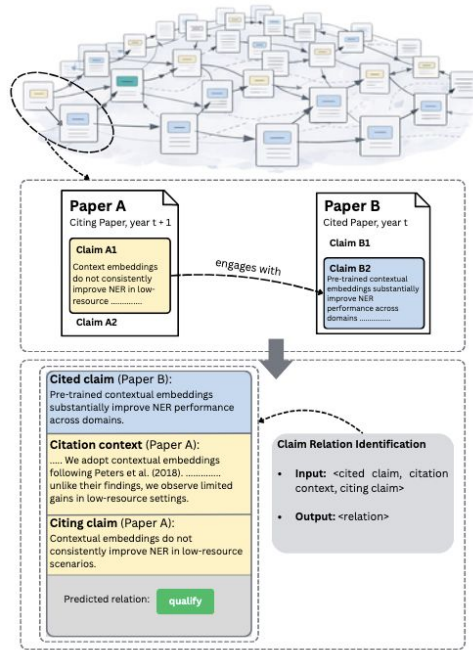
- **Cross-document scientific claim relations classification grounded in citations:** analyze whether a citing paper **supports**, **extends**, **qualifies**, **refutes**, or **references** a claim as background [1]



Claims are most often reused as contextual premises or explicitly supported by subsequent work.

Citation Analysis: ClaimFlow

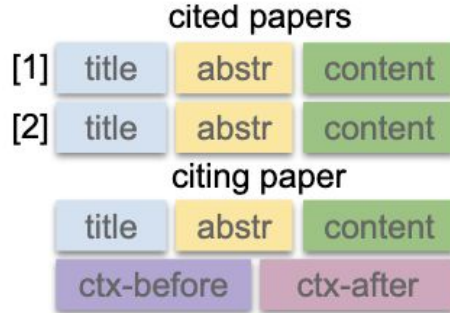
- **Cross-document scientific claim relations classification grounded in citations:** analyze whether a citing paper **supports**, **extends**, **qualifies**, **refutes**, or **references** a claim as background [1]



only a small number of claims are ever challenged (11.1%) at least once

Related Work and Citation Text Generation

- Generating **citation text** for pre-selected cited papers given **the context** of the citing paper

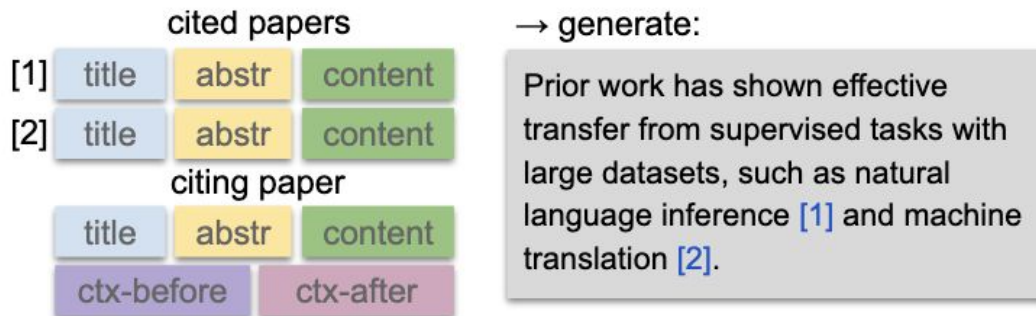


→ generate:

Prior work has shown effective transfer from supervised tasks with large datasets, such as natural language inference [1] and machine translation [2].

Related Work and Citation Text Generation

- Generating **citation text** for pre-selected cited papers given **the context** of the citing paper



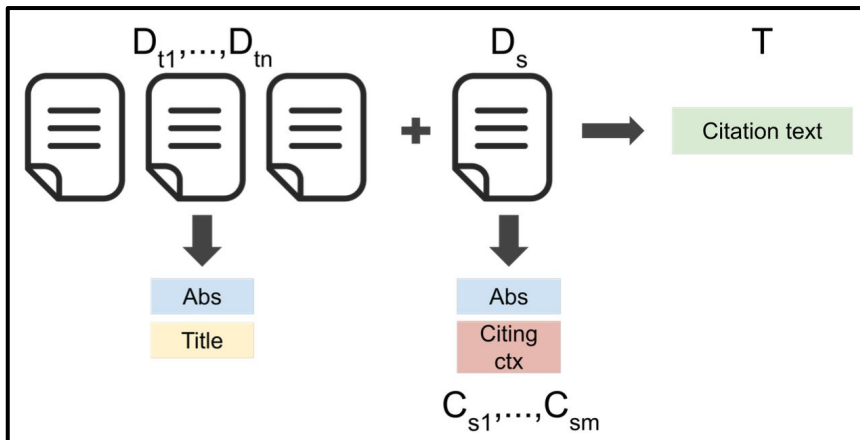
- Task variations

Dataset	Input					Output		
	Cited document (D^t)			Citing context (C^s)		Citation text (T)		
	Single	Abs	Multi Abs	Title	Abs	Text	Sent	Para
AbuRa'ed et al. (2020)	✓			✓			✓	
Chen et al. (2021)			✓					✓
Lu et al. (2020)			✓		✓			✓
Xing et al. (2020)	✓				✓	✓	✓	

Related work section generation

CiteBench

- A benchmark for citation text generation^[1]
 - **Unified task definition and dataset**
 - **Standardized baselines:** unsupervised, supervised, transfer-based
 - **Evaluation kit:** standard metrics and discourse-based measurements

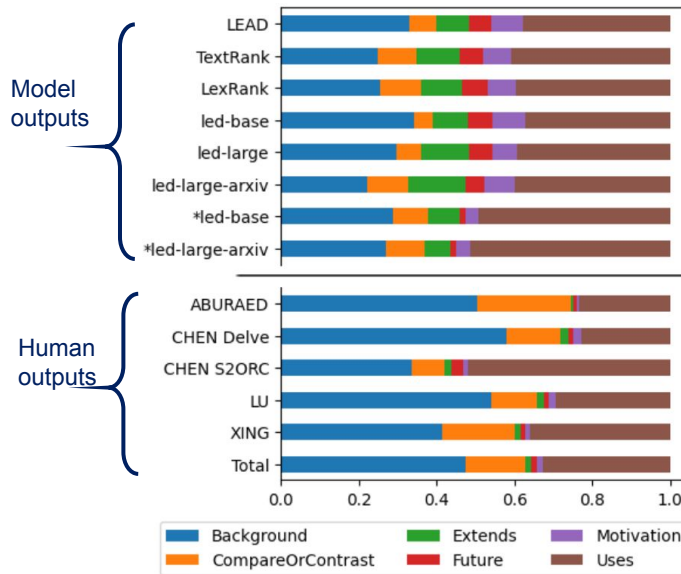
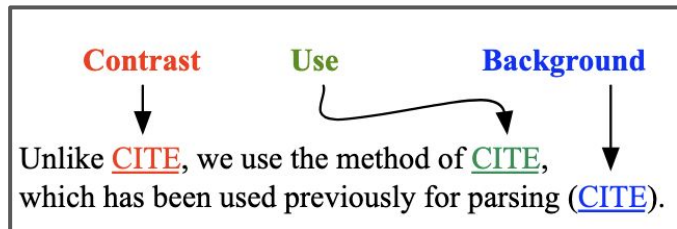


[1] Funkquist et al., EMNLP 2023. [CiteBench: A Benchmark for Scientific Citation Text Generation](#).

CiteBench

- Discourse analysis based on **citation intent**^[1] and CORWA (citation oriented related work annotation)

citation intent analysis example

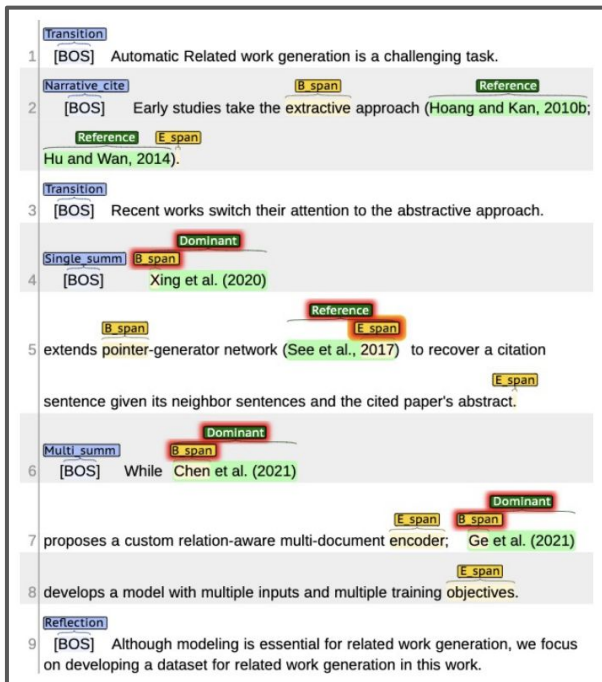


Models tend to under-generate the *Background* and *CompareOrContrast* sentences, they produce more *Future*, *Uses* and *Extends* sentences than the gold reference.

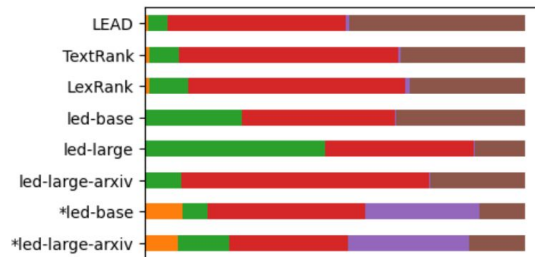
[1] Jurgen et al., TACL 2018. [Measuring the Evolution of a Scientific Field through Citation Frames.](#)

- Discourse analysis based on citation intent and **CORWA** (citation oriented related work annotation)^[1]

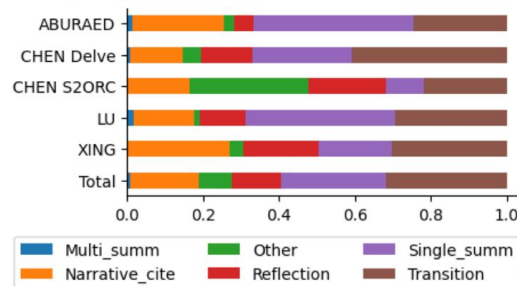
CORWA analysis example



Model outputs

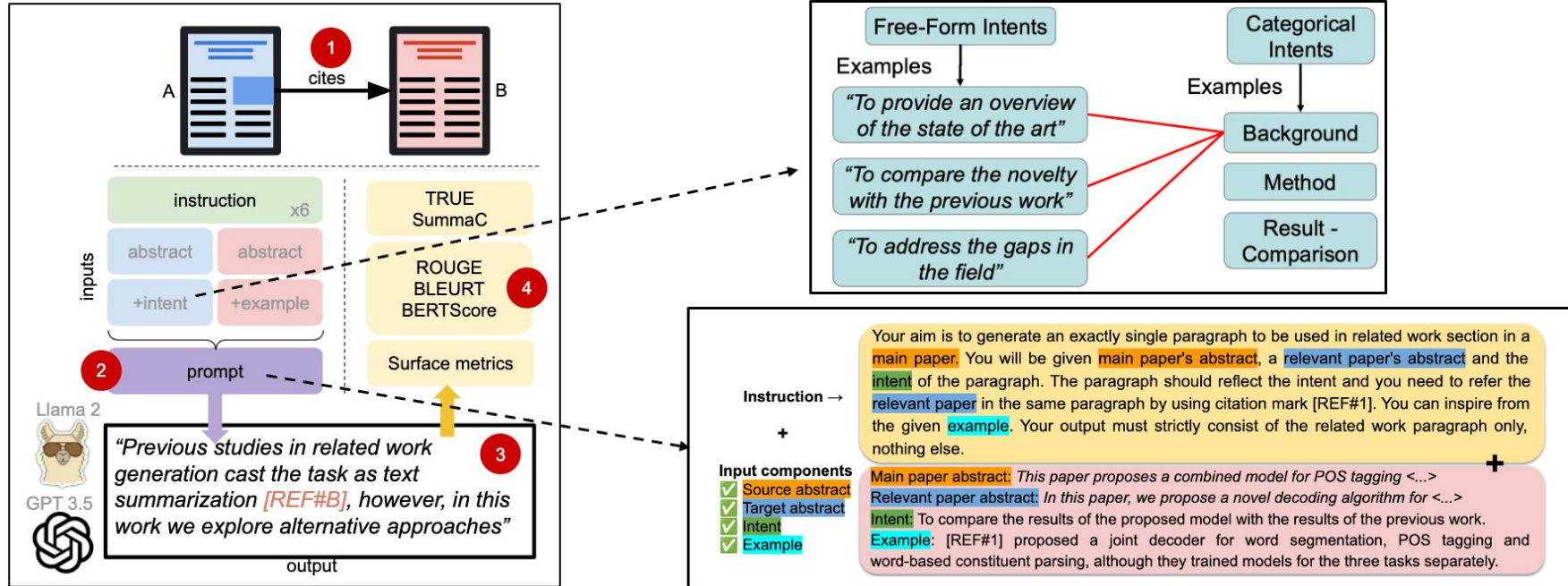


Human outputs



Models tend to under-generate the **Narrative_cite** and **Single_summ** class, while over-generating the **Reflection**, compared to the distributions in the gold reference texts.

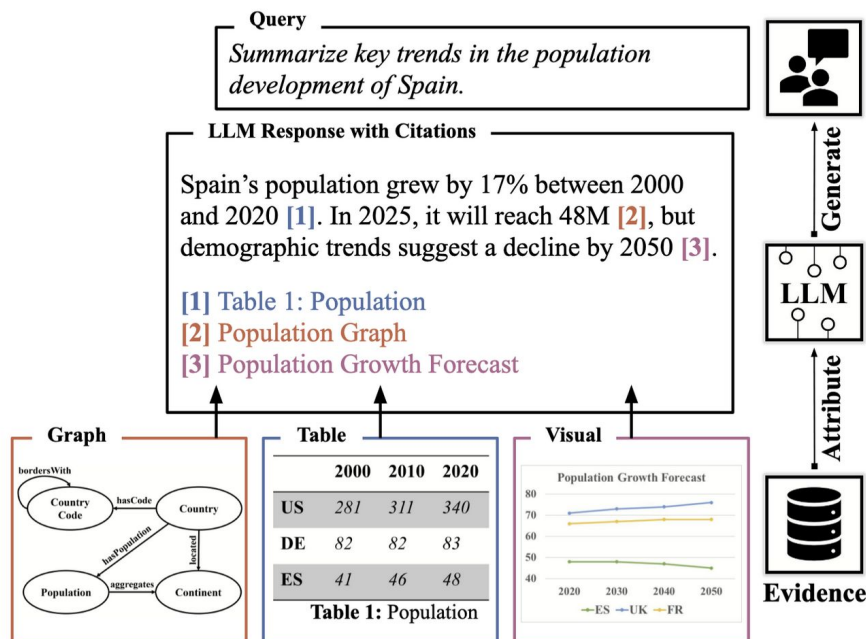
Key Important Input Components for Citation Text Generation



- Joint use of **citation intent** and **example sentences** gives best results for both models and human annotators
- **Free-form** citation intents are more effective than **categorical** intents

LLM Citation Generation and Attribution

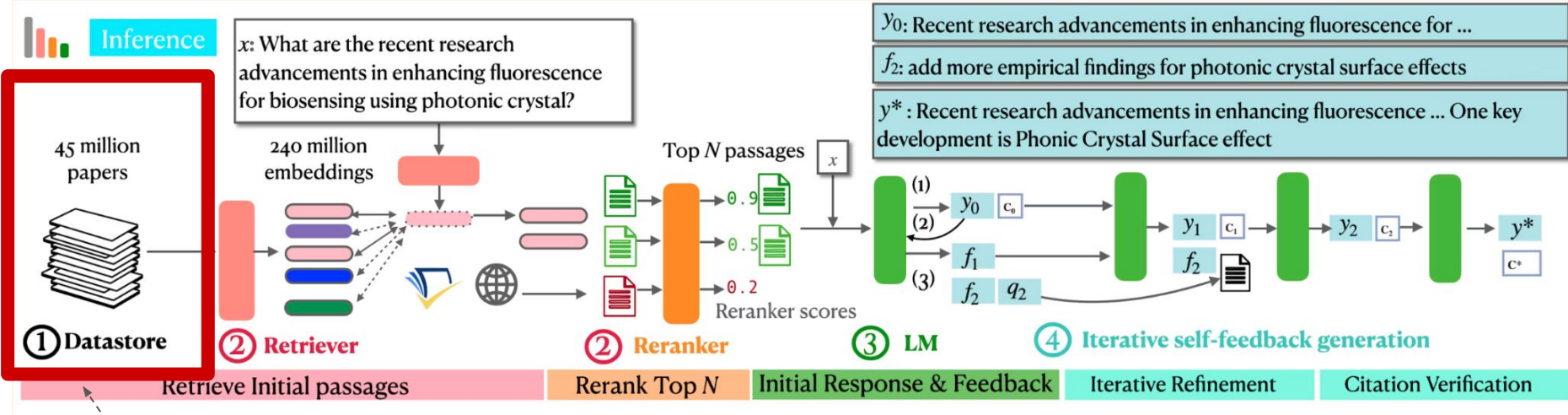
- **Evidence-based text generation**: given a query, LLMs generate a response where each claim is grounded with inline citations to identifiable source passages, making outputs verifiable.



[1] Li et al., arXiv 2023. [A Survey of Large Language Models Attribution](#).

[2] Schreieder et al., arXiv 2025. [Attribution, Citation, and Quotation: A Survey of Evidence-based Text Generation with Large Language Models](#).

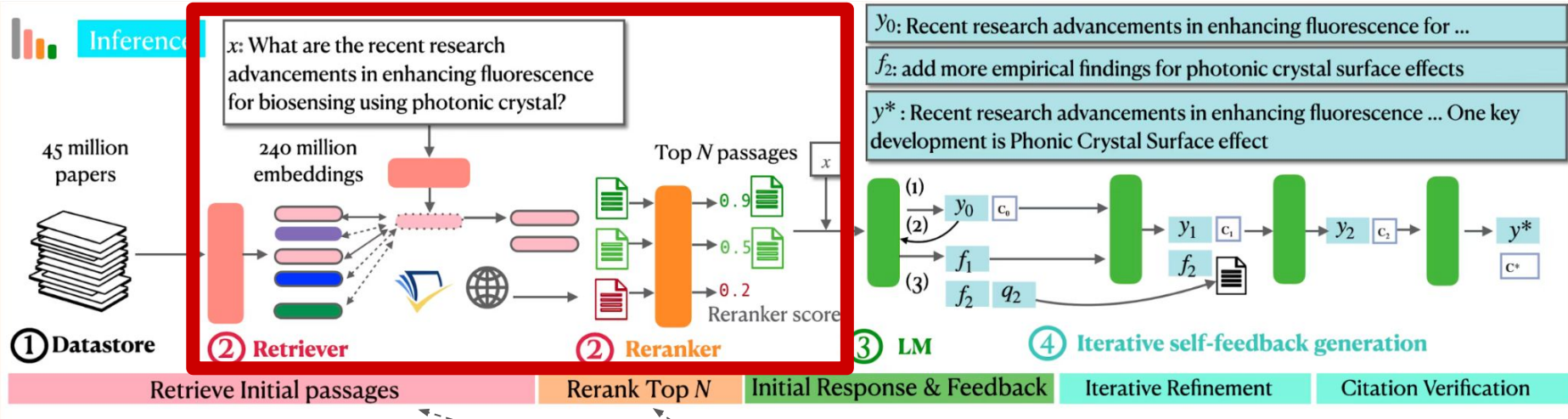
- Inference



A collection of more than 45M papers from Semantic Scholar and ~250M corresponding passage embeddings.

[1] Asai et al., Nature 2026. [Synthesizing scientific literature with retrieval-augmented language models.](#)

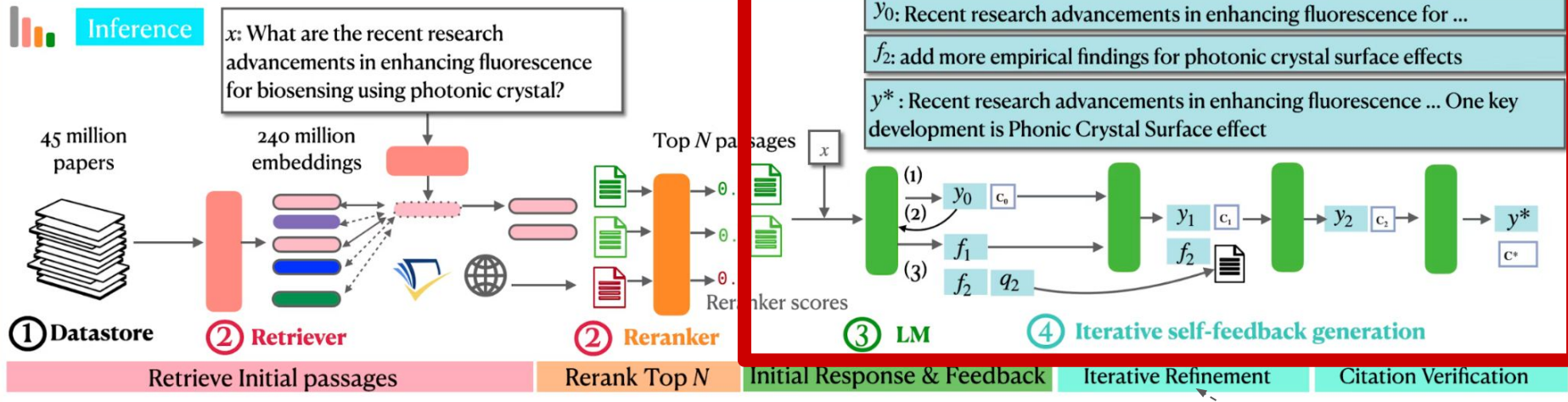
- Inference



Specialized Retrievers and Rerankers:
These tools are trained specifically to identify relevant passages from our scientific literature datastore

[1] Asai et al., Nature 2026. [Synthesizing scientific literature with retrieval-augmented language models.](#)

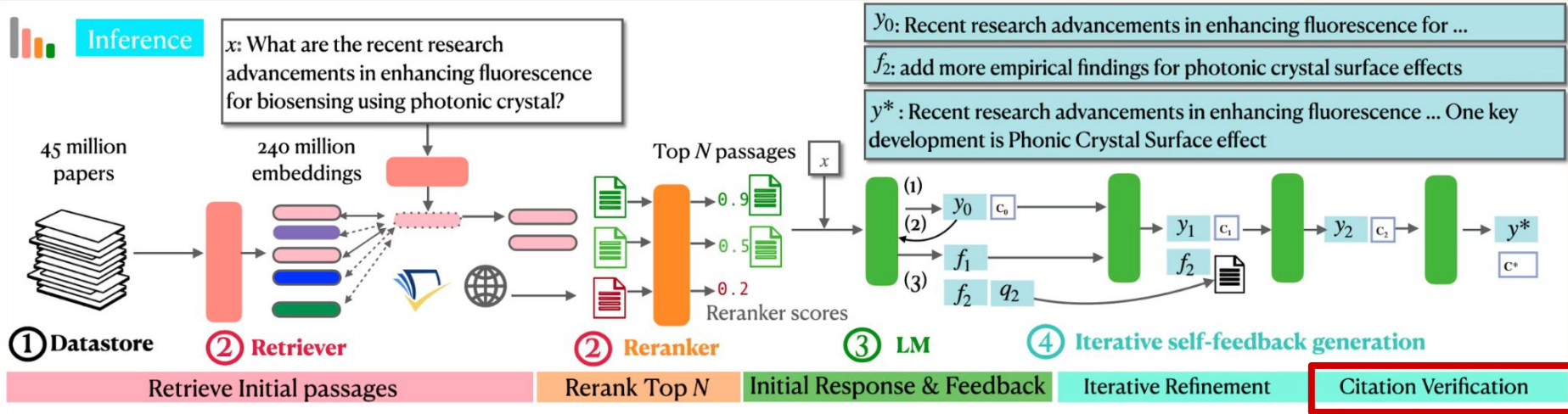
- Inference



Iterative self-feedback to refine model outputs through natural language feedback. Each iteration involves additionally retrieving more papers.

[1] Asai et al., Nature 2026. [Synthesizing scientific literature with retrieval-augmented language models.](#)

- Inference

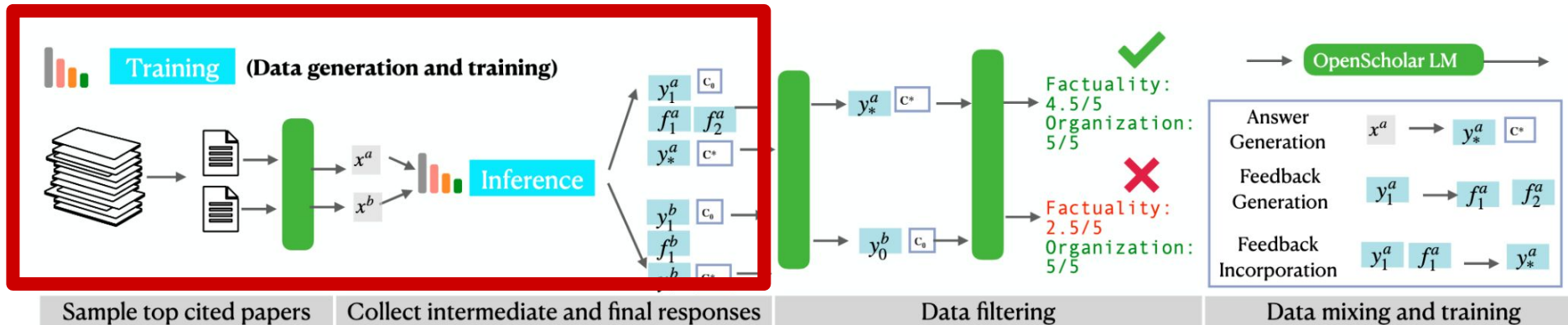


The generator LM ensures that all citation-worthy statements are adequately supported by references from the retrieved passages. If any claims lack proper citations, the LM performs a **post hoc insertion** to ensure that citation-worthy statements are supported by passages.

[1] Asai et al., Nature 2026. [Synthesizing scientific literature with retrieval-augmented language models.](#)

OpenScholar

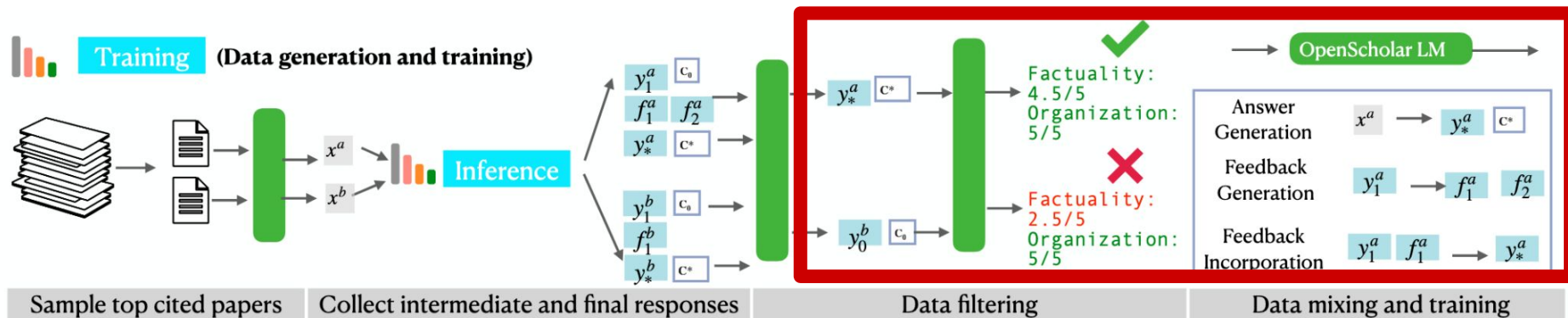
- Training LLama 3.1 8B



LLama 3.1 70B

generating information-seeking queries based on their abstracts that may require multiple papers to answer

- Training LLama 3.1 8B



LLama 3.1 70B

- 1) pair-wise filtering, compare the quality of model outputs at the the final step and the initial step, and retain the output that is judged to be higher quality at the final step.
- 2) Rate the selected response on a 5-point scale for two criteria: **organization**, and **factual precision/citation accuracy**. To be valid, the model's output must score at least 4.5 in both.

- ScholarQA-CS
 - 100 questions + detailed answer rubrics
 - Expert annotators (expert annotators holding Ph.D.s in the field)

Input

What are the best practices to protect a software against vulnerabilities from third party libraries?

Rubric

Must Have Item-1: The answer should discuss best practices that can be used to prevent these implications, such as reliable source, update monitoring, code

Nice to Have Item-1: The answer could provide some examples of famous third-party libraries that can be used in different programming languages.

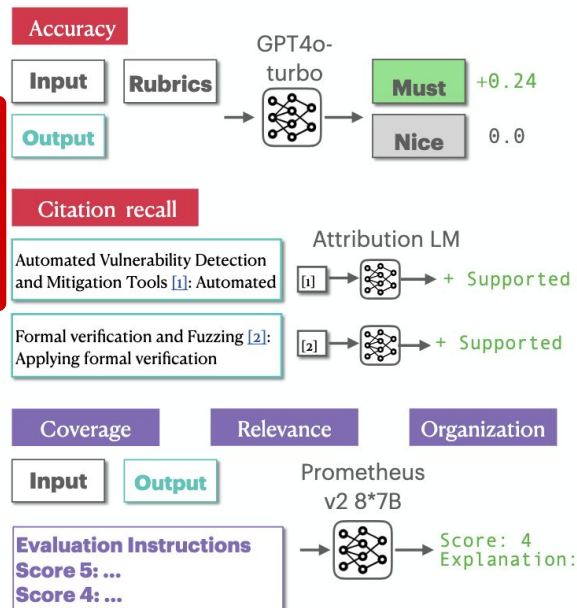
Output Must have item-1 is included ✓

Protecting software against vulnerabilities stemming from third-party libraries is a crucial aspect of software security [1] [2]. Below are some of the best practices based on the existing literature:

Citations

[1] To solve the challenges faced by third-party libraries, researchers can take the following measures: (i) Develop intelligent security tools to automatically detect and repair vulnerabilities in third-party libraries ...

[2] Applying formal verification methods to examine the security properties of ... can also be part of library test suites or continuous integration in order to run



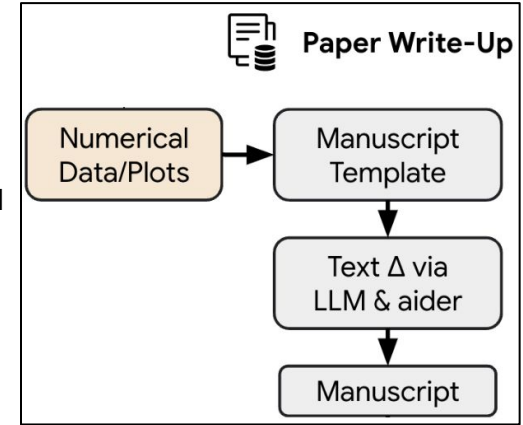
Agenda

- ❑ **Short Text Generation with Citation Grounding**
- ❑ **Automatic Research Paper Writing**
- ❑ **Automated Survey & Deep Research Generation**
- ❑ **Meta-analysis Table Generation and Comparative Literature Synthesis**
- ❑ **Summary**

The AI Scientist Paper Writing-up Process

Step 1: Per-section text generation

1. Prompt Aider to fill in a blank conference latex paper template section by section
2. All previous sections of the paper has already written are in the context of the language model
3. Each section is initially refined with one round of self-reflection
4. Aider is prompted to not include any citations in the text at this stage, and fill in only a skeleton for the related work



Paper Writing Aider Prompt

We've provided the `latex/template.tex` file to the project. We will be filling it in section by section.

First, please fill in the `{section}` section of the writeup.

Some tips are provided below:

`{per_section_tips}`

Before every paragraph, please include a brief description of what you plan to write in that paragraph in a comment.

Be sure to first name the file and use `*SEARCH/REPLACE*` blocks to perform these edits.

How to ML Paper - A brief Guide

Feel free to [comment / share](#) and happy paper writing! Also, please see caveats* below. If you like this, why not follow [How to ML](#) on Twitter and share the advice/love?

Canonical ML Paper Structure

Abstract (TL;DR of paper):

X: What are we trying to do and why is it relevant?

Y: Why is this hard?

Z: How do we solve it (i.e. our contribution)?

1: How do we verify that we solved it:

1a) Experiments and results

1b) Theory

Introduction (Longer version of the **Abstract**, i.e. of the entire paper):

X: What are we trying to do and why is it relevant?

Y: Why is this hard?

Z: How do we solve it (i.e. our contribution)?

1: How do we verify that we solved it:

1a) Experiments and results, including comparison to prior SOTA if applicable

1b) Theory

2: New trend: specifically list your contributions as bullet points (credits to [Brandon](#))

Extra space? Future work!

Extra points for having [Figure 1](#) on the first page

Related Work:

Academic siblings of our work, i.e. alternative attempts in literature at trying to solve the same problem.

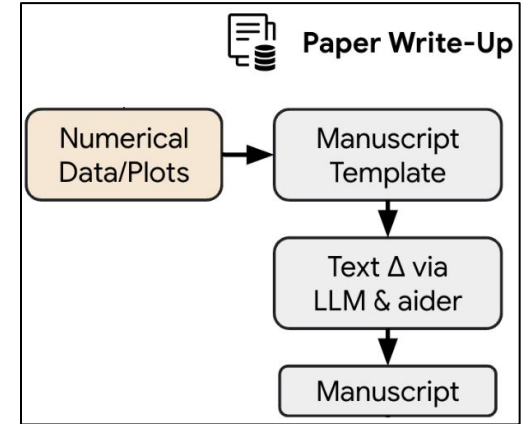
Goal is to "Compare and contrast" - how does their approach differ in either assumptions or method? If their method is applicable to our **Problem Setting** I expect a comparison in the experimental section. If not, there needs to be a clear statement why a given method is not applicable.

Note: Just describing what another paper is doing is not enough. We need to compare and contrast.

The AI Scientist Paper Writing-up Process

Step 2: Web search for references

1. 20 rounds to poll the Semantic Scholar API looking for the most relevant sources to compare and contrast the near-completed paper against for the **related work section**
2. Select papers for discussion and complete missing citations elsewhere in the paper
3. For each selected paper, a **brief note on where/how to cite it is sent to Aider**, and its BibTeX is automatically added to the LaTeX file



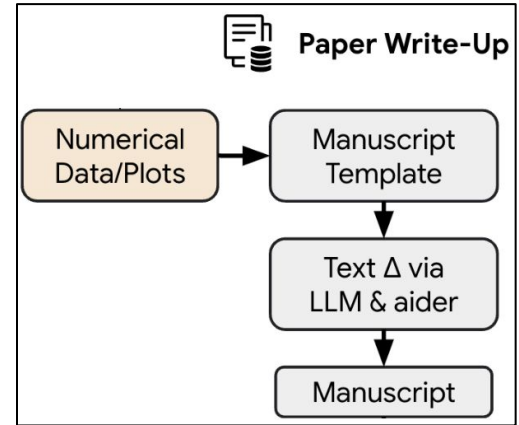
The AI Scientist Paper Writing-up Process

Step 2: Web search for references

1. 20 rounds to poll the Semantic Scholar API looking for the most relevant sources to compare and contrast the near-completed paper against for the **related work section**
2. Select papers for discussion and complete missing citations elsewhere in the paper
3. For each selected paper, a **brief note on where/how to cite it is sent to Aider**, and its BibTeX is automatically added to the LaTeX file

Step 3: Refinement

One final round of self-reflection section-by-section (aiming to remove any duplicated information and streamline the arguments)



The AI Scientist Paper Writing-up Process

Step 2: Web search for references

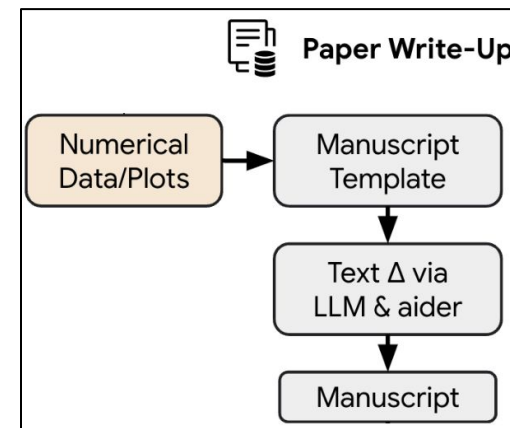
1. 20 rounds to poll the Semantic Scholar API looking for the most relevant sources to compare and contrast the near-completed paper against for the **related work section**
2. Select papers for discussion and complete missing citations elsewhere in the paper
3. For each selected paper, a **brief note on where/how to cite it is sent to Aider**, and its BibTeX is automatically added to the LaTeX file

Step 3: Refinement

One final round of self-reflection section-by-section (aiming to remove any duplicated information and streamline the arguments)

Step 4: Compilation

After the Latex paper template is populated, it is compiled, and any errors are sent back to Aider for automatic correction



Agent Skills: ML-Paper-Writing

- An AI agent skill that transforms a research repo with code & experimental results into a publication-ready LaTeX paper, targeting top ML/AI and systems venues.

- Workflow

- **Understand:** Explore the repo, identify contribution, search literature
- **Draft:** Write proactively section by section, with the scientist giving feedback at each step
- **Cite:** Never hallucinate citations, always fetch BibTeX programmatically via Semantic Scholar/DOI; mark unverified as [placeholder]
- **Format & Submit:** Use conference LaTeX templates, check page limits, run checklists



Orchestra-Research / AI-Research-SKILLS

Workflow 0: Starting from a Research Repository

When beginning paper writing, start by understanding the project:

Project Understanding:

- [] Step 1: Explore the repository structure
- [] Step 2: Read README, existing docs, and key results
- [] Step 3: Identify the main contribution with the scientist
- [] Step 4: Find papers already cited in the codebase
- [] Step 5: Search for additional relevant literature
- [] Step 6: Outline the paper structure together
- [] Step 7: Draft sections iteratively with feedback

Step 1: Explore the Repository

```
# Understand project structure
ls -la
find . -name "*.py" | head -20
find . -name "*.md" -o -name "*.txt" | xargs grep -l -i "result|conclusion|finding"
```

Look for:

- README.md - Project overview and claims
- results/, outputs/, experiments/ - Key findings
- configs/ - Experimental settings
- Existing .bib files or citation references
- Any draft documents or notes

Step 2: Identify Existing Citations

Check for papers already referenced in the codebase:

```
# Find existing citations
grep -r "arxiv|doi|cite" --include="*.md" --include="*.bib" --include="*.py"
find . -name "*.bib"
```

These are high-signal starting points for Related Work—the scientist has already deemed them relevant.

Step 3: Clarify the Contribution

Before writing, explicitly confirm with the scientist:

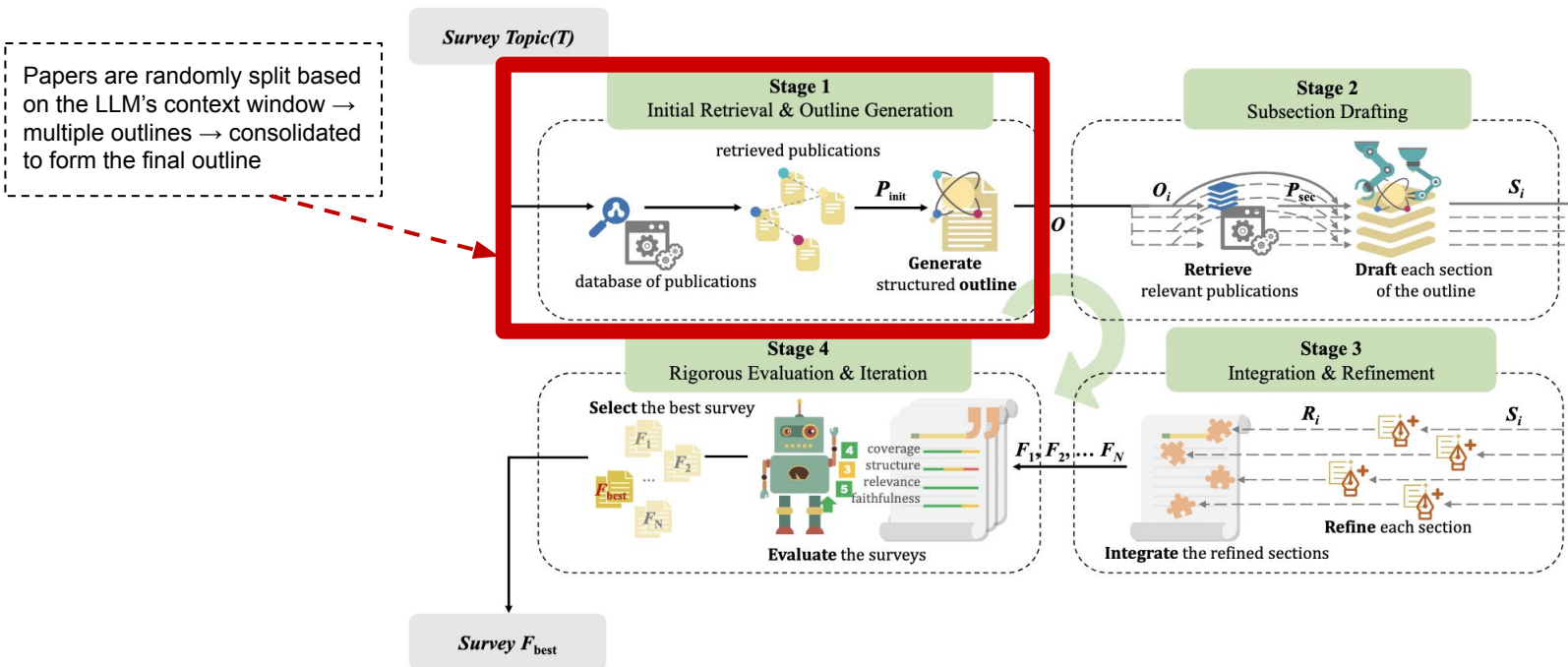
"Based on my understanding of the repo, the main contribution appears to be [X]. The key results show [Y]. Is this the framing you want for the paper, or should we emphasize different aspects?"

Agenda

- ❑ **Short Text Generation with Citation Grounding**
- ❑ **Automatic Research Paper Writing**
- ❑ **Automated Survey & Deep Research Generation**
- ❑ **Meta-analysis Table Generation and Comparative Literature Synthesis**
- ❑ **Summary**

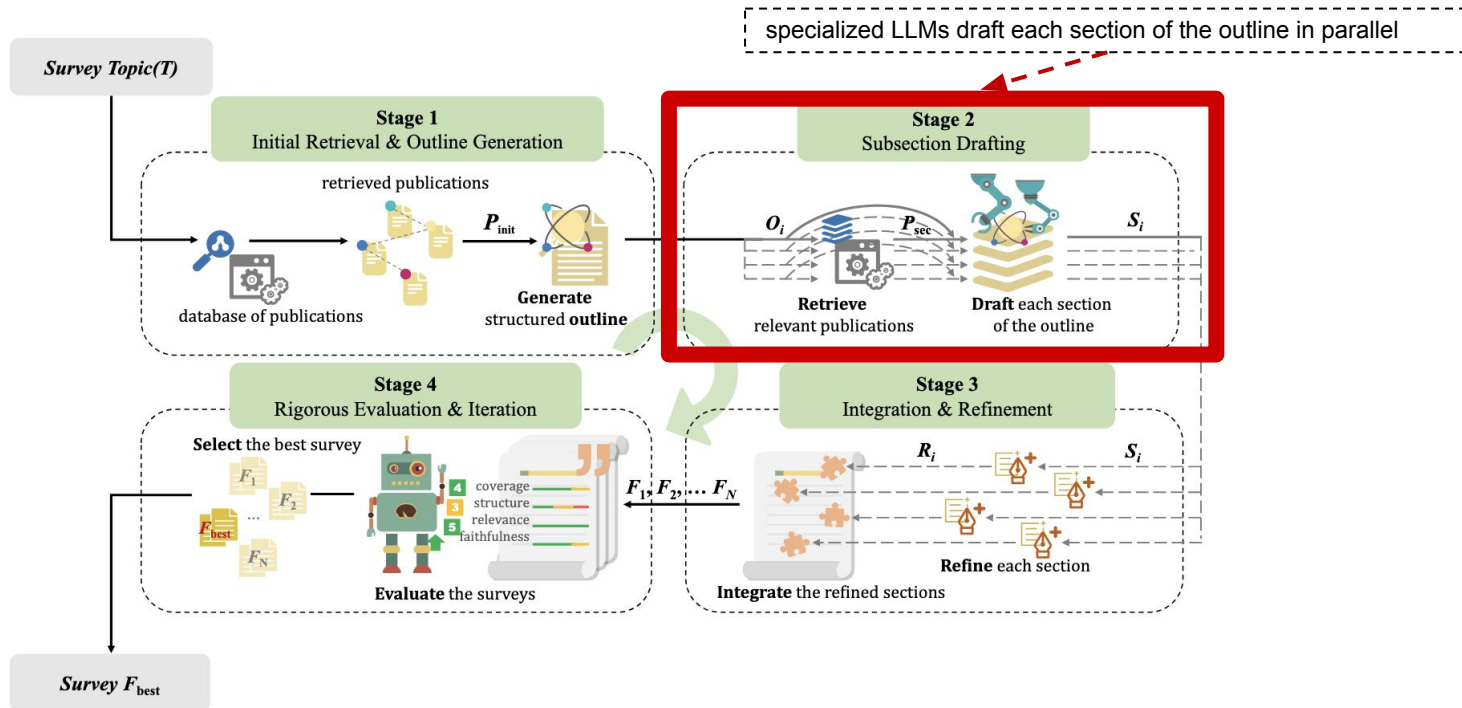
AutoSurvey

- AutoSurvey pipeline - a cost of \$1.2 and 3 minutes per survey (with Claude-haiku)



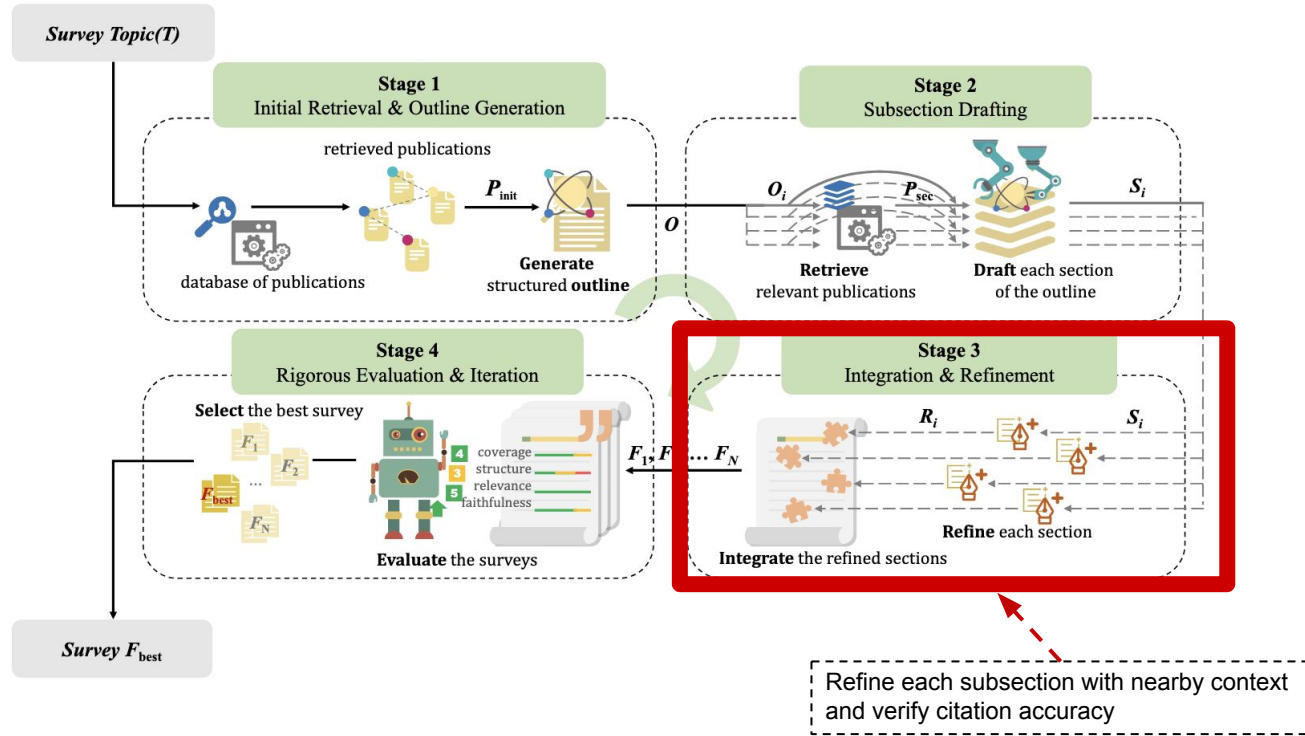
AutoSurvey

- AutoSurvey pipeline - a cost of \$1.2 and 3 minutes per survey (with Claude-haiku)



AutoSurvey

- AutoSurvey pipeline - a cost of \$1.2 and 3 minutes per survey (with Claude-haiku)



- Evaluation on 20 survey papers about LLMs (Claude-haiku as the writer)

Survey Title	Citations
A survey for in-context learning	323
A Survey on Large Language Models for Recommendation	55
A Survey of Detecting LLM-Generated Texts	42
Explainability for Large Language Models	25
A Survey on Evaluation of Large Language Models	183
A Survey on Large Language Model based Autonomous Agents	101
A Survey of Large Language Models in Medicine	234
Domain Specialization as the Key to Make Large Language Models Disruptive	14
Practical and Ethical Challenges of Large Language Models in Education	53
Aligning Large Language Models with Human	53
A Survey on ChatGPT and Beyond	144
Instruction Tuning for Large Language Models	45
Large Language Models for Information Retrieval	22
Towards Safer Generative Language Models: Safety Risks, Evaluations, and Improvements	17
A Survey of Chain of Thought Reasoning	13
A Survey on Hallucination in Large Language Models	116
Bias and Fairness in Large Language Models	12
Large-scale Multi-Modal Pre-trained Models	61
A Survey on Model Compression and Acceleration for Pretrained Language Models	22
Large Language Models for Software Engineering	49

- Evaluation on 20 survey papers about LLMs (Claude-haiku as the writer)

Survey Length (#tokens)	Methods	Speed	Citation quality		Coverage	Content Quality		Avg.
			Recall	Precision		Structure	Relevance	
8k	Human writing	0.16	80.00	87.50	4.50	4.16	5.00	4.52
	Naive RAG-based LLM generation	79.67	78.14 \pm 5.23	71.92 \pm 6.83	4.40 \pm 0.48	3.86 \pm 0.71	4.86 \pm 0.33	4.33
	AutoSurvey	107.00	82.48 \pm 2.77	77.42 \pm 3.28	4.60 \pm 0.48	4.46 \pm 0.49	4.8 \pm 0.39	4.61
16k	Human writing	0.14	88.52	79.63	4.66	4.38	5.00	4.66
	Naive RAG-based LLM generation	43.41	71.48 \pm 12.50	65.31 \pm 15.36	4.46 \pm 0.49	3.66 \pm 0.69	4.73 \pm 0.44	4.23
	AutoSurvey	95.51	81.34 \pm 3.65	76.94 \pm 1.93	4.66 \pm 0.47	4.33 \pm 0.59	4.86 \pm 0.33	4.60
32k	Human writing	0.10	88.57	77.14	4.66	4.50	5.00	4.71
	Naive RAG-based LLM generation	22.64	79.88 \pm 4.35	65.03 \pm 8.39	4.41 \pm 0.64	3.75 \pm 0.72	4.66 \pm 0.47	4.23
	AutoSurvey	91.46	83.14 \pm 2.44	78.04 \pm 3.14	4.73 \pm 0.44	4.26 \pm 0.69	4.8 \pm 0.54	4.58
64k	Human writing	0.07	86.33	77.78	5.00	4.66	5.00	4.88
	Naive RAG-based LLM generation	12.56	68.79 \pm 11.00	61.97 \pm 13.45	4.4 \pm 0.61	3.66 \pm 0.47	4.66 \pm 0.47	4.19
	AutoSurvey	73.59	82.29 \pm 3.64	77.41 \pm 3.84	4.73 \pm 0.44	4.33 \pm 0.47	4.86 \pm 0.33	4.62

LLM as a judge

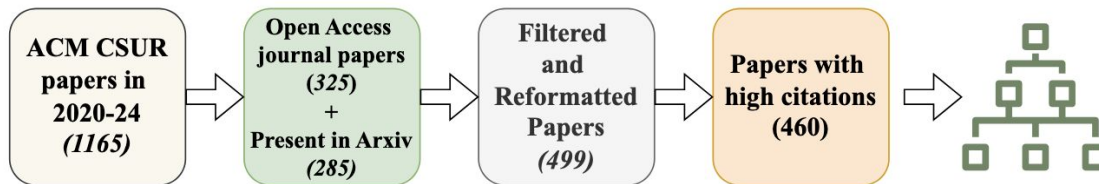
73.59 survey papers per hour

NLI model to assess the claim grounding coverage and relevance of citations

"It achieves near-human levels of coverage, relevance, and citation quality while maintaining a significantly lower time cost."

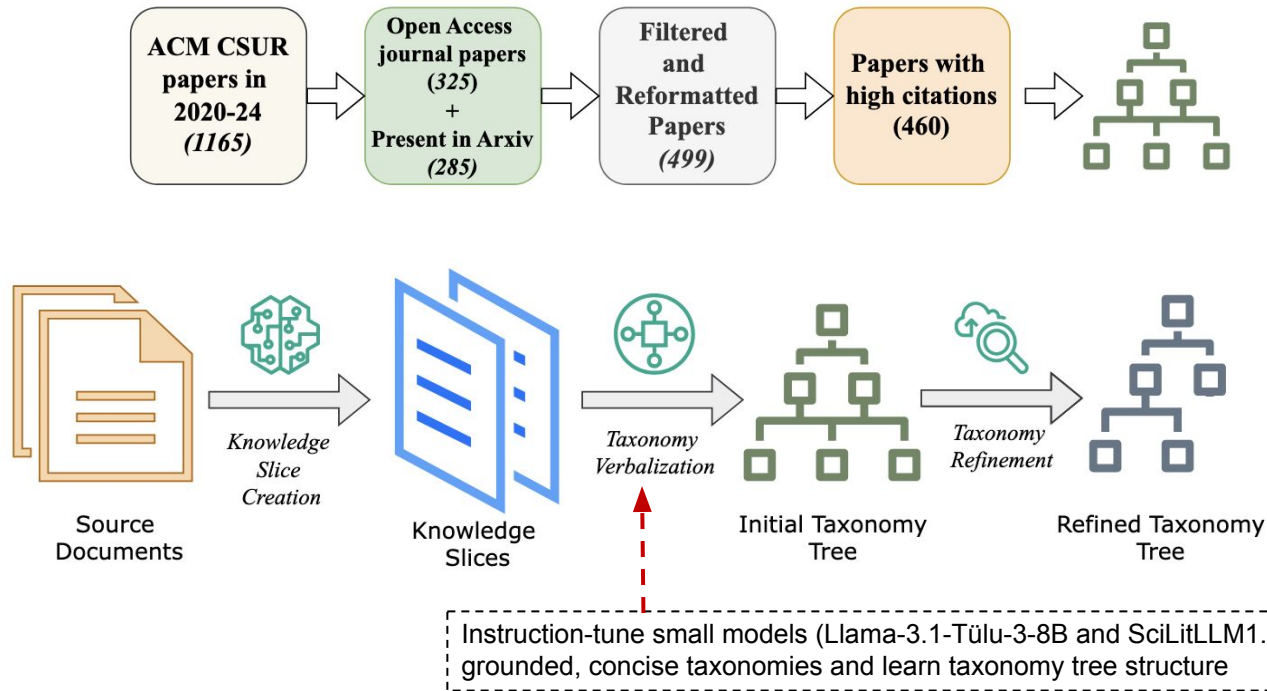
TaxoAlign

- Goal: Automated taxonomy creation that can bridge the gap between human-generated and automatically-created taxonomies



TaxoAlign

- Goal: Automated taxonomy creation that can bridge the gap between human-generated and automatically-created taxonomies



Agenda

- ❑ **Short Text Generation with Citation Grounding**
- ❑ **Automatic Research Paper Writing**
- ❑ **Automated Survey & Deep Research Generation**
- ❑ **Meta-analysis Table Generation and Comparative Literature Synthesis**
- ❑ **Summary**

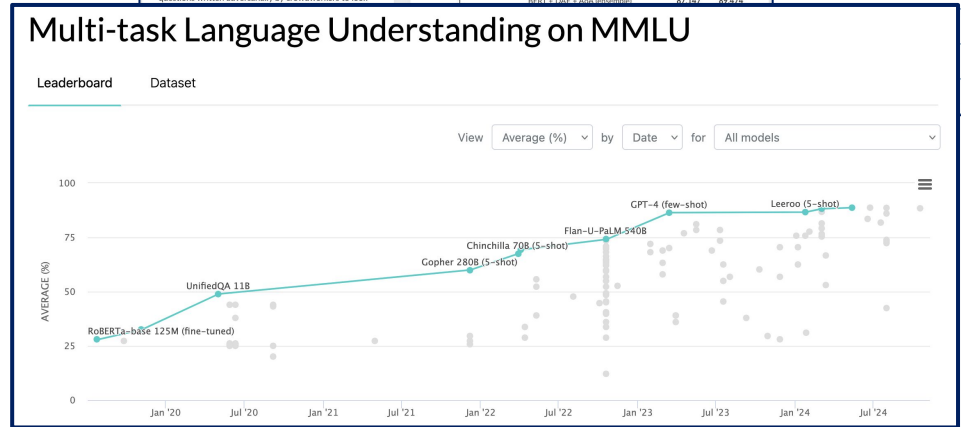
Meta-analysis for NLP/ML Literature: Scientific Leaderboard Construction

- A **scientific leaderboard** is a ranked list of methods/models evaluated on the same task using a shared benchmark and metric.

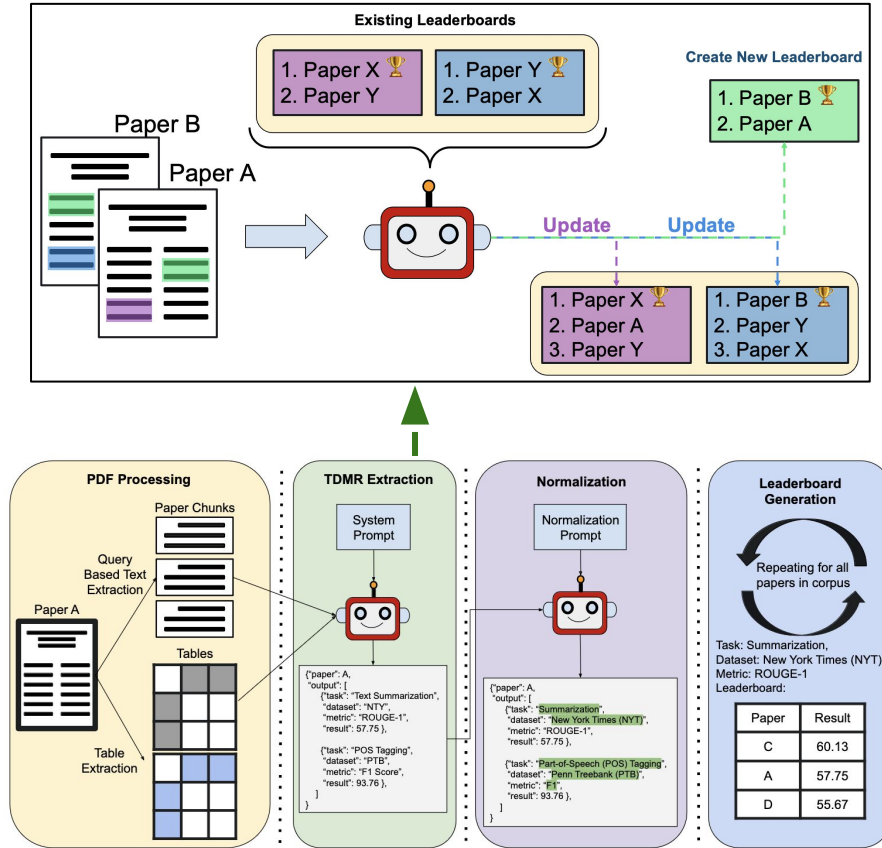


Automate the Whole Process

Rank	Model	EM	F1	QQP	MNLI-m	MNLI-mm	QNLI	RTE	WNLI	AX
1	Human Performance Stanford University (Rajpurkar & Ji et al. 18)	86.831	89.452	74.2/90.3	90.2	89.8	98.6	86.3	90.4	47.5
2	BERT + DAE (Kosaraju et al. 2020)	87.147	89.674	73.7/89.9	87.9	87.4	96.0	86.3	89.0	42.8
3				59.5/80.4	92.0	92.8	91.2	93.6	95.9	NaN
4				74.4/90.7	88.2	87.9	95.7	83.5	80.8	43.9
5				73.1/89.9	87.6	87.2	93.9	80.9	65.1	39.9
6				73.2/89.8	89.1	88.5	94.0	76.0	71.9	44.7
7				84.5	65.1	42.4				
8				80.4	65.1	40.7				
9				79.8	65.1	28.3				



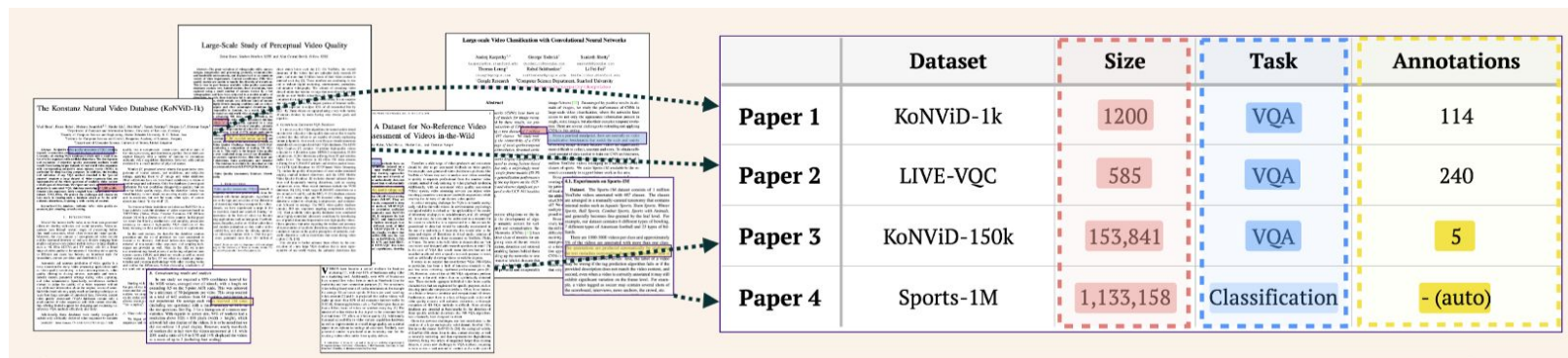
Meta-analysis for NLP/ML Literature: Scientific Leaderboard Construction



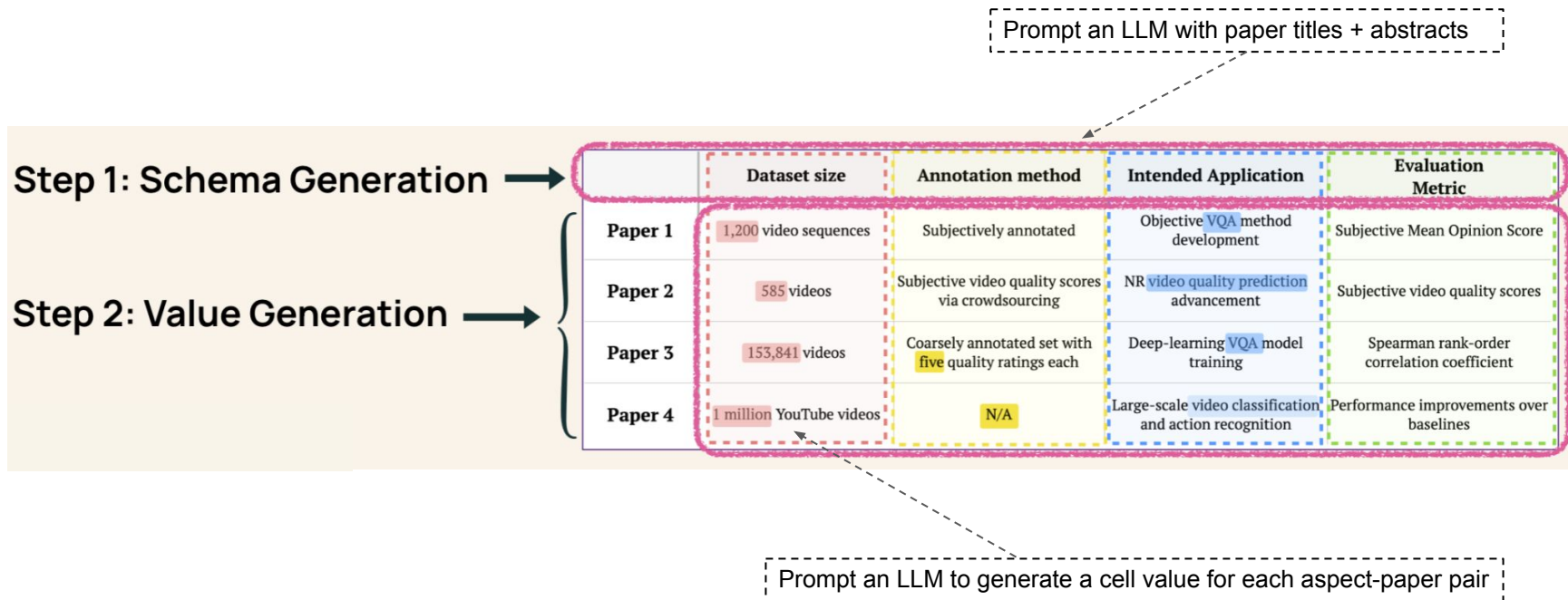
ArxivDIGESTables: Automatically Generate Literature Review Tables

- **Task: literature review table generation**

- Rows are a set of papers
- Columns are a set of aspects that the papers share



- Two stage decomposition generation

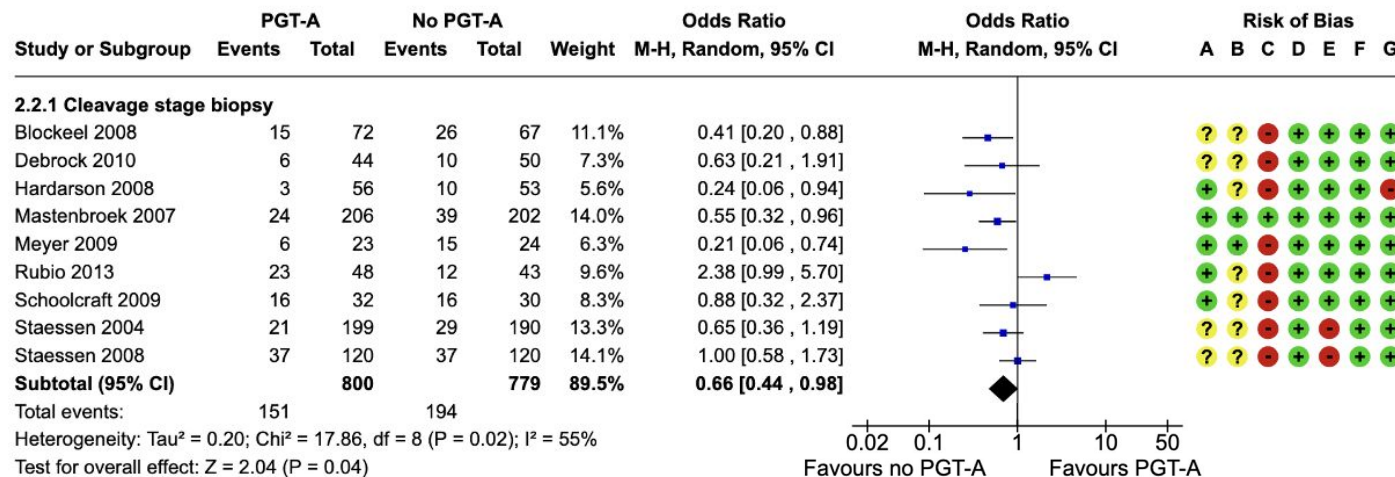


Meta-analysis for Biomedical Literature: Forest-plot Generation

- A forest plot is a chart that shows the results of several studies in one place.
- It helps you compare each study and see the overall result when all studies are combined.

Meta-analysis for Biomedical Literature: Forest-plot Generation

Research question: Whether the intervention PGT-A (formerly known as preimplantation genetic screening) leads to fewer live births than in the control group for women undergoing an IVF (in vitro fertilization) treatment?

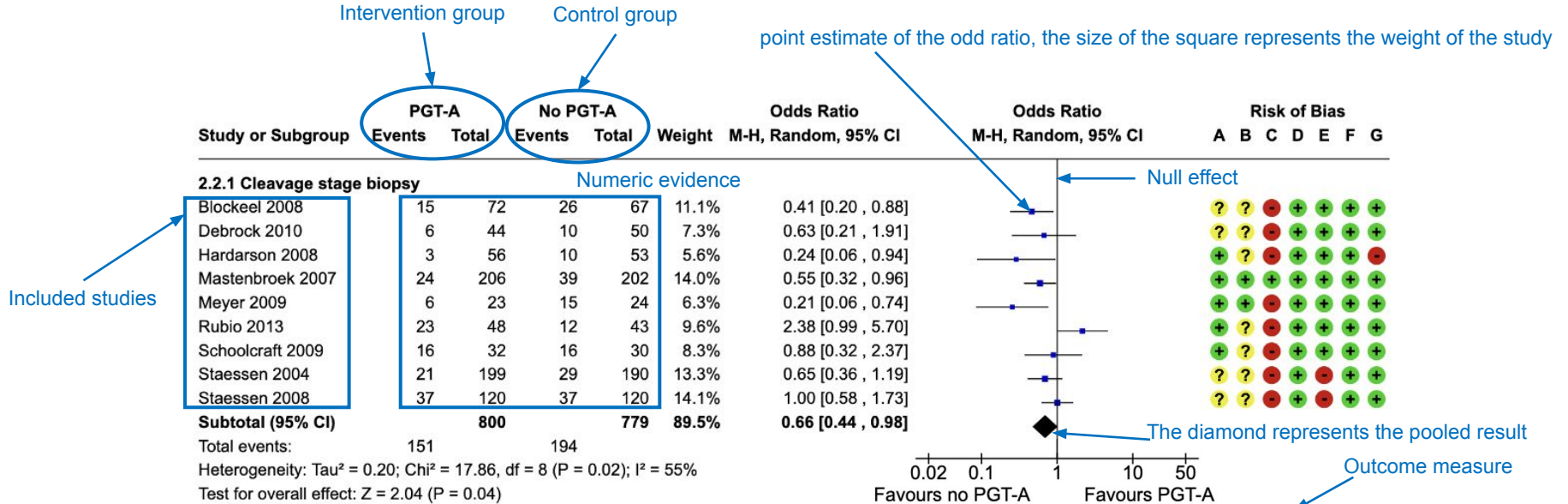


IVF with PGT-A versus IVF without PGT-A with the use of FISH for the genetic analysis, outcome: live birth rate after the first embryo transfer per woman randomised.

Conclusions: The currently available evidence is insufficient to support PGT-A in routine clinical practice.

Meta-analysis for Biomedical Literature: Forest-plot Generation

Research question: Whether the intervention PGT-A (formerly known as preimplantation genetic screening) leads to fewer live births than in the control group for women undergoing an IVF (in vitro fertilization) treatment?

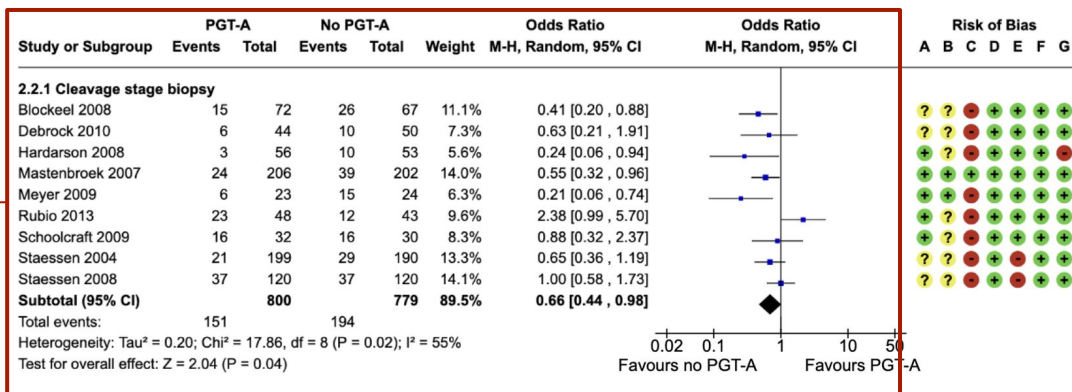
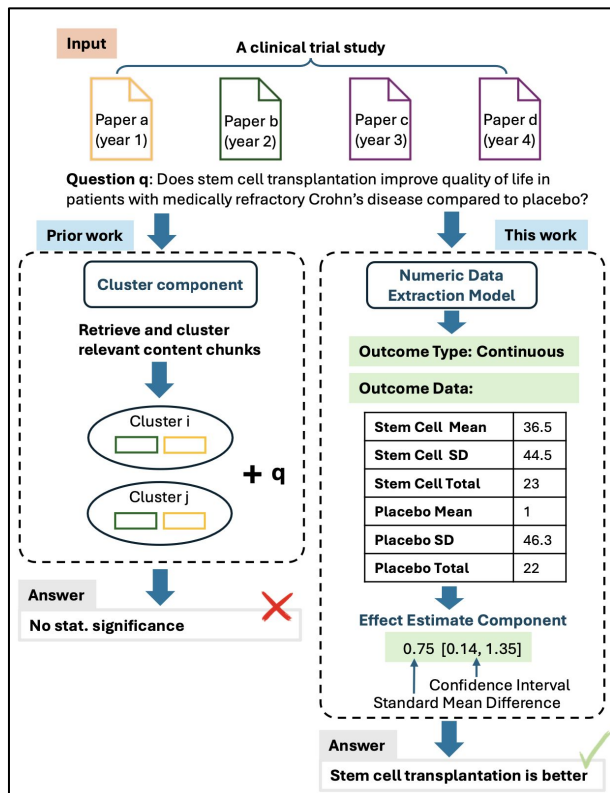


IVF with PGT-A versus IVF without PGT-A with the use of FISH for the genetic analysis, outcome: live birth rate after the first embryo transfer per woman randomised.

Conclusions: The currently available evidence is insufficient to support PGT-A in routine clinical practice.

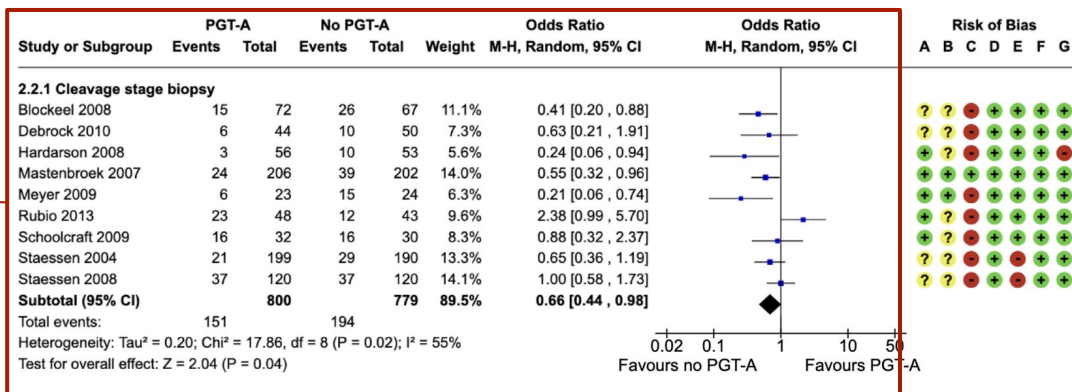
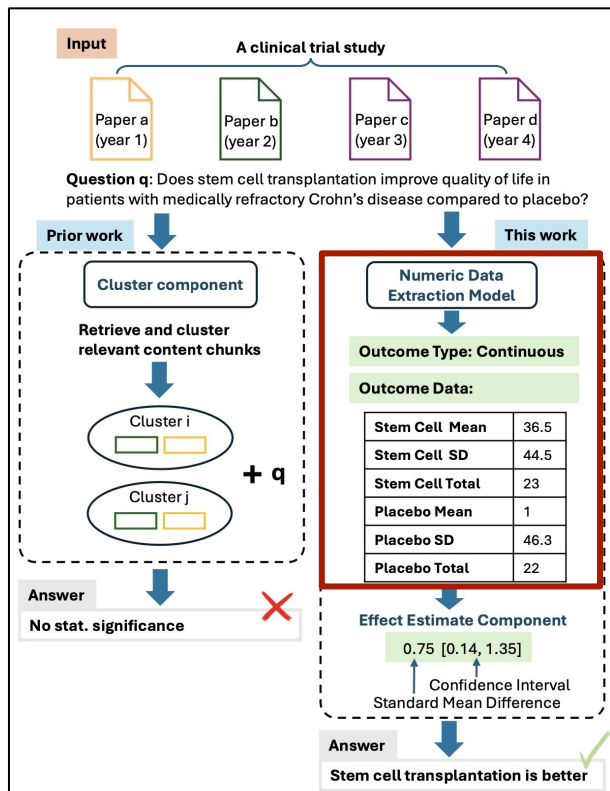
Meta-analysis for Biomedical Literature: Forest-plot Generation

- A numeric reasoning approach to estimate study effect



Meta-analysis for Biomedical Literature: Forest-plot Generation

- A numeric reasoning approach to estimate study effect



Fine-tune a numeric evidence extraction model using RL with GRPO

- Correctness Reward (CR):

$$R_{CR} = \frac{1 + \sum_{j=1}^n 1\{v_j \approx \hat{v}_j\}}{1 + n}$$

- Format Reward (FR):

$$R_{FR} = \begin{cases} 1 & \text{if } \pi_{\theta}(x) \in \mathcal{F} \\ 0 & \text{otherwise} \end{cases}, \text{ with } \mathcal{F} \text{ set of valid formats}$$

- Thought Format Reward (TFR):

$$R_{TFR} = \begin{cases} 1 & \text{if } \pi_{\theta}(x) \text{ matches thought pattern} \\ 0 & \text{otherwise} \end{cases}$$

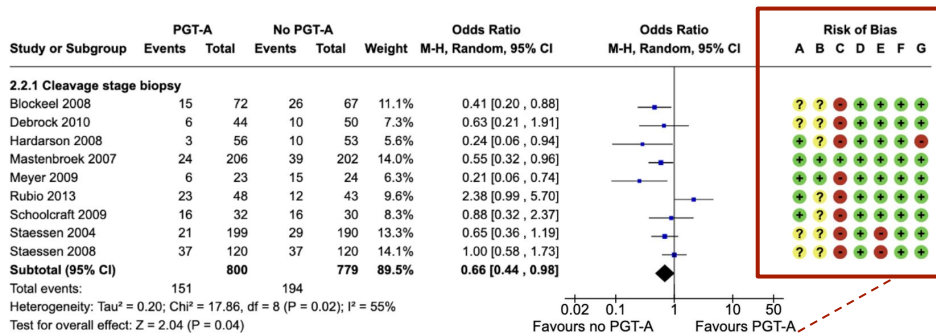
- Final reward:

$$R = 0.8 \cdot R_{CR} + 0.1 \cdot R_{FR} + 0.1 \cdot R_{TFR}$$

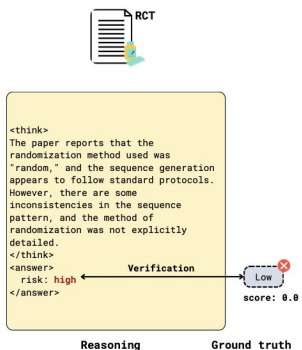
Existing reasoning models are insufficient for complex domain-specific tasks.

Meta-analysis for Biomedical Literature: Forest-plot Generation

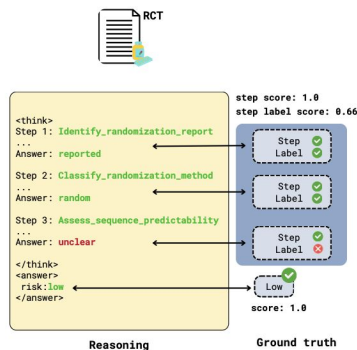
- Verifiable process reward models for risk of bias analysis



Verifiable Outcome Rewarding



Verifiable Process Rewarding



$$R(Y; x) = \sum_{t=1}^T r_t(Y; x) + r_{\text{label}}$$

(Process) (Outcome)

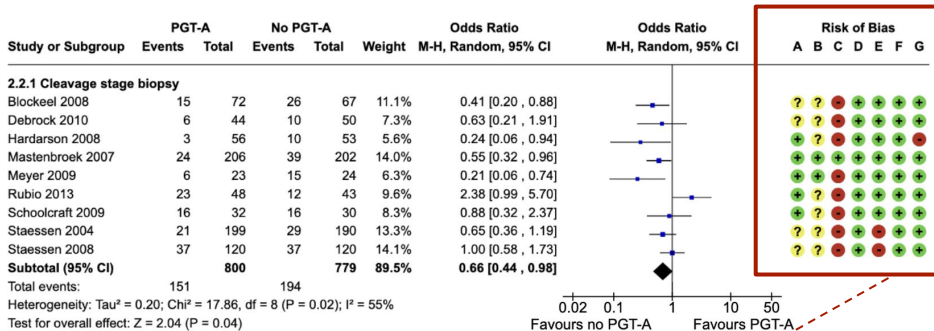
$$r_t(Y; x) = w_t^0 s_t^0(s_t, s_t^*) + w_t^1 s_t^1(\hat{\ell}_t, \ell_t^*)$$

(Step name) (Step label)

The process reward is a weighted sum of correct step identification and correct label selection.

Meta-analysis for Biomedical Literature: Forest-plot Generation

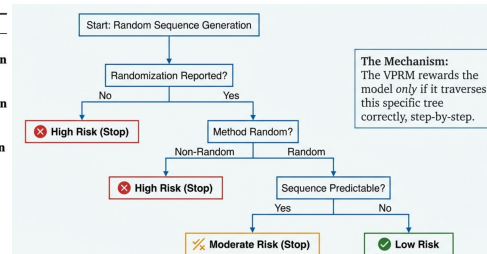
- Verifiable process reward models for risk of bias analysis



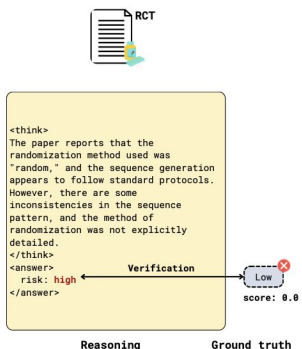
Error caught immediately during the process. Model penalized.

```

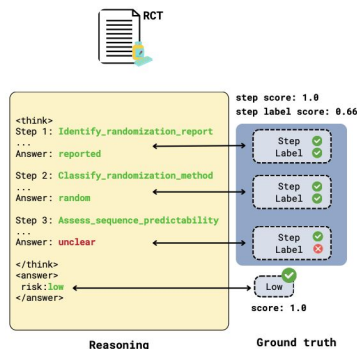
Algorithm 1 RoB A Macro
1: procedure PREDICTLABEL-A(steps)
2:   if steps[IDENTIFYRANDOMIZATIONREPORT] = NOTREPORTED then
3:     return MODERATE
4:   end if
5:   if steps[CLASSIFYRANDOMIZATIONMETHOD] = NONRANDOM then
6:     return HIGH
7:   end if
8:   if steps[ASSESSSEQUENCEPREDICTABILITY] = PREDICTABLE then
9:     return MODERATE
10:  end if
11:  if steps[BASELINEIMBALANCE] = LIKELY then
12:    return HIGH
13:  end if
14:  return LOW
15: end procedure
    
```



Verifiable Outcome Rewarding



Verifiable Process Rewarding



$$R(Y; x) = \sum_{t=1}^T r_t(Y; x) + r_{label}$$

(Process) (Outcome)

$$r_t(Y; x) = w_t^p s_t^p(s_t, s_t^*) + w_t^l s_t^l(\hat{\ell}_t, \ell_t^*)$$

(Step name) (Step label)

The process reward is a weighted sum of correct step identification and correct label selection.

Agenda

- ❑ **Short Text Generation with Citation Grounding**
- ❑ **Automatic Research Paper Writing**
- ❑ **Automated Survey & Deep Research Generation**
- ❑ **Meta-analysis Table Generation and Comparative Literature Synthesis**
- ❑ **Summary**

	Pre-LLM Era	LLM Era
Key tasks	Related work generation Citation text generation Cite-worthiness detection Citation recommendation Citation intent classification Citation analysis	LLM citation generation and attribution Agentic deep research <ul style="list-style-type: none">● Planning● Information gathering● Tool usage● Synthesis and reporting
Core challenges	Coherent multi-doc synthesis	Verifiable, non-hallucinated attribution
Generation quality	Often disfluent or extractive	Fluent but may hallucinate or post-rationalize
Evaluation	Rouge, F1, BERTScore	LLM-as-a-judge, citation quality based on NLI models (citation precision/citation recall)

- **The verification bottleneck:** Generation has outpaced verification!
 - The field can now produce fluent, cited text at scale
 - Automated attribution verification (NLI-based, LLM-as-judge) is still unreliable that human checking remains necessary.

- **The verification bottleneck:** Generation has outpaced verification!



- **AI-assisted scientific writing**
 - Fluent AI-generated text with seemingly grounded citations can create a false sense of rigor.
 - Without reliable automated verification, AI paper-writing tools may accelerate scientific output while **quietly eroding the trust and reliability on which researchers depend.**

References

1. Wright & Augenstein, ACL 2021 Findings. [CiteWorth: Cite-Worthiness Detection for Improved Scientific Document Understanding.](#)
2. Çelik & Tekir, EMNLP 2025. [CiteBART: Learning to Generate Citations for Local Citation Recommendation.](#)
3. Jurgen et al., TACL 2018. [Measuring the Evolution of a Scientific Field through Citation Frames.](#)
4. Li et al., NAACL 2022. [CORWA: A Citation-Oriented Related Work Annotation Dataset.](#)
5. Pramanick et al., arXiv 2026. [ClaimFlow: Tracing the Evolution of Scientific Claims in NLP.](#)
6. Chen et al., ACL 2021. [Capturing Relations between Scientific Papers: An Abstractive Model for Related Work Section Generation.](#)
7. Lu et al., EMNLP 2020. [Multi-XScience: A Large-scale Dataset for Extreme Multi-document Summarization of Scientific Articles.](#)
8. Xing et al., ACL 2020. [Automatic Generation of Citation Texts in Scholarly Papers: A Pilot Study.](#)
9. Funkquist et al., EMNLP 2023. [CiteBench: A Benchmark for Scientific Citation Text Generation.](#)
10. Şahinuç et al., ACL 2024. [Systematic Task Exploration with LLMs: A Study in Citation Text Generation.](#)
11. Li et al., arXiv 2023. [A Survey of Large Language Models Attribution.](#)
12. Schreieder et al., arXiv 2025. [Attribution, Citation, and Quotation: A Survey of Evidence-based Text Generation with Large Language Models.](#)
13. Asai et al., Nature 2026. [Synthesizing scientific literature with retrieval-augmented language models.](#)
14. Lu et al., arXiv 2024. [The AI Scientist: Towards Fully Automated Open-Ended Scientific Discovery.](#)
15. Wang et al., NeurIPS 2024. [AutoSurvey: Large Language Models Can Automatically Write Surveys.](#)
16. Lahiri et al., EMNLP 2025. [TaxoAlign: Scholarly Taxonomy Generation Using Language Models.](#)
17. Hou et al., ACL 2019. [Identification of Tasks, Datasets, Evaluation Metrics, and Numeric Scores for Scientific Leaderboards Construction.](#)
18. Şahinuç et al., EMNLP 2024. [Efficient Performance Tracking: Leveraging Large Language Models for Automated Construction of Scientific Leaderboards.](#)
19. Newman et al., EMNLP 2024. [ArxivDIGESTables: Synthesizing Scientific Literature into Tables using Language Models.](#)
20. Pronesti et al., EMNLP 2025. [Enhancing Study-Level Inference from Clinical Trial Papers via RL-based Numeric Reasoning.](#)
21. Pronesti et al., arXiv 2026. [Beyond Outcome Verification: Verifiable Process Reward Models for Structured Reasoning.](#)