
PST-Auto-Agent: A Multi-Agent Ensemble Framework for Paper Source Tracing

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The escalating volume of scientific literature necessitates efficient methods for
2 identifying foundational works that significantly inform new research. This paper
3 addresses the Paper Source Tracing (PST) problem, which aims to quantify the
4 influence of cited references on a focal paper, assigning importance weights to
5 its most salient sources. To this end, we propose a novel multi-agent ensemble
6 architecture for PST, integrating Deepseek-R1-250528, GPT-5-2025-08-07, and
7 Gemini-2.5-pro. Our system employs a robust pipeline, featuring advanced XML
8 parsing, empirically optimized prompt engineering with counterfactual reasoning
9 and multi-role Socratic dialogue, and a sophisticated multi-agent integration strat-
10 egy. This strategy utilizes weighted model predictions, intelligent default scoring,
11 and a consistency penalty mechanism to derive precise source paper identifications.
12 Our method becomes a strong tuning-free baseline for the PST problem that does
13 not require feature engineering. Our method also achieves top-ranked results when
14 combined with feature engineering techniques. This work highlights the efficacy
15 of multi-agent ensembles and advanced prompt engineering for complex academic
16 information tracing tasks.

1 Introduction

18 The proliferation of academic literature across various scientific domains necessitates efficient
19 methods for understanding the intellectual lineage and foundational influences of research papers.
20 Identifying the primary sources that significantly inform a paper’s central ideas or fundamental
21 methodologies is crucial for accurate academic attribution, knowledge graph construction, and the
22 broader understanding of scientific evolution. This challenge defines the *Paper Source Tracing (PST)*
23 *problem*: given a focal paper P and its full text, the goal is to identify its most salient references, here
24 termed *source papers*, and quantify their influence with an importance weight. A reference qualifies
25 as a primary source if paper p ’s central idea or fundamental methodology is rooted in it.

26 Traditional approaches to reference analysis often rely on simple citation counts or textual similarity,
27 which frequently fall short in discerning the true intellectual contributions of cited works. Some
28 advanced methods finetune large language models (LLMs) or graph neural networks (GNNs) by
29 employing long texts and citation structures (Zhong et al., 2024). However, these methods are
30 resource-intensive. Many cited papers serve merely as background information, datasets, software
31 tools, or benchmarking studies, rather than embodying the core ideas or methods of the focal paper.
32 The inherent complexity of distinguishing between merely cited and genuinely influential sources
33 underscores the need for sophisticated, domain-aware analytical tools. This problem is further
34 compounded by the sheer volume of publications and the intricate web of interconnections within
35 scientific literature.

36 To address the limitations of existing methods and provide a robust benchmark for the PST problem,
37 we introduce **PST-Bench** (Zhang et al., 2024), a novel dataset comprising 2,141 meticulously labeled

38 computer science publications. This dataset was constructed through a rigorous annotation process
39 involving computer science graduate students specializing in relevant subfields, who identified source
40 papers within their respective areas of expertise. The collaborative online paper group workflow,
41 coupled with extensive expert validation and preprocessing, ensures the high quality and reliability of
42 PST-Bench. PST-Bench serves as an invaluable resource for developing and evaluating PST solutions,
43 providing both professionally annotated data and supplementary rule-generated annotations.

44 In this paper, we propose a novel multi-agent ensemble architecture designed to tackle the PST
45 problem effectively. Our system leverages the combined strengths of state-of-the-art large language
46 models (LLMs): Deepseek-R1-250528, GPT-5-2025-08-07, and Gemini-2.5-pro. The system oper-
47 ates through a structured pipeline: XML Preprocessing, Prompt Engineering, Multi-Agent Prediction,
48 Intelligent Ensemble, and Prediction Method. This architecture is meticulously designed for ro-
49 bustness and scalability, facilitating the precise identification of foundational methodological and
50 conceptual papers within a citation network.

51 **XML Processing and Data Extraction** is critical for handling diverse XML structures, employing
52 a dual-parsing strategy with primary and fallback parsers, complemented by a comprehensive data
53 cleaning pipeline to ensure 100% reliability. **Prompt Engineering Strategy** is grounded in extensive
54 empirical analysis, evaluating over 1,000 human-annotated papers to identify optimal configurations.
55 Our unified prompt architecture incorporates advanced reasoning frameworks such as Counterfactual
56 Reasoning, Idea DNA Matching, and Multi-Role Socratic Dialogue, alongside strict exclusion criteria
57 to filter non-source references.

58 The core of our approach lies in the **Multi-Agent Integration Strategy**, which employs an advanced
59 ensemble methodology. This intelligently combines predictions from the three complementary
60 LLMs with optimized weight allocation. The ensemble algorithm also designs Consistency Penalty
61 Mechanism to mitigate inconsistent predictions. The final prediction score for each LLM is a product
62 of its penalty score, weight score, and initial prediction score.

63 We evaluate our proposed method using Mean Average Precision (MAP), a standard metric for
64 ranking tasks, calculated as the average of Average Precision (AP) across all papers in the test set.
65 Our experimental results on PST-Bench demonstrate that the *pst-auto-agent* model significantly
66 outperforms individual baseline models. Our ensemble model, *pst-auto-agent*, achieved the
67 highest performance with a MAP score of 0.388, representing a notable 22.0% relative improvement
68 over the best baseline Gemini-2.5-pro. This superior performance validates the architectural design
69 and the effectiveness of our multi-agent ensemble approach for the PST task.

70 Furthermore, we showcase the practical utility of our tuning-free method by integrating it into the top-
71 ranked solution English Hercules in KDD Cup 2024. English Hercules is a GPU-free approach that
72 combines feature engineering and LLM API-based methods. By ensembling GPT-5, DeepSeek-R1,
73 and Gemini-pro into its framework using the ensemble method of English Hercules, we observed a
74 clear complementary effect, significantly enhancing the overall performance. The improved ranking
75 on the KDD Cup 2024 leaderboard underscores the robustness and broad applicability of our multi-
76 agent ensemble strategy.

77 In summary, this paper makes the following key contributions:

- 78 • We formally define the Paper Source Tracing problem and introduce **PST-Bench** (Zhang et al.,
79 2024), a new, expertly annotated dataset comprising 2,141 computer science publications with
80 rigorous quality control and temporal partitioning for robust evaluation.
- 81 • We propose a novel multi-agent ensemble architecture, *pst-auto-agent*, which integrates
82 multiple state-of-the-art LLMs with advanced prompt engineering and intelligent ensemble
83 strategies.
- 84 • We demonstrate the superior performance of our proposed method on PST-Bench, achieving a
85 MAP score of 0.388, significantly outperforming strong baseline models.
- 86 • We illustrate the practical impact of our approach by showing its ability to enhance the perfor-
87 mance of a leading method in the KDD Cup 2024 competition.

88 These contributions pave the way for more accurate and automated identification of foundational
89 intellectual contributions in scientific literature, thereby enriching academic research and knowledge
90 discovery.

91 **2 Related Work**

92 Paper source tracing builds upon several related research areas, including citation analysis, bibliometrics,
93 and scientific literature mining. Early work in citation analysis focused on identifying influential
94 papers through citation counts and network centrality measures Chubin & Garfield (1980). More
95 recent approaches have incorporated machine learning techniques to identify seminal works and
96 research trends.

97 Bibliometric studies have explored various aspects of scientific communication, including co-citation
98 analysis Small (1973), bibliographic coupling KESSLER (1963), and topic modeling Blei et al.
99 (2003). These methods provide valuable insights into the structure of scientific knowledge but often
100 lack the precision needed for accurate source tracing.

101 Recent advances in natural language processing and graph neural networks have enabled more
102 sophisticated approaches to literature analysis. Methods such as document embedding Le & Mikolov
103 (2014), graph convolutional networks Kipf & Welling (2017), and transformer models Vaswani et al.
104 (2017) have been applied to scientific text analysis with promising results. And paper source tracing
105 have witnessed significant methodological innovations, particularly within the context of the KDD
106 Cup 2024 OAG Challenge. The top-performing solutions introduced novel architectures and learning
107 paradigms that extend beyond traditional citation analysis and bibliometric modeling.

108 Chen et al. (2024a) graft BERT predictions from noisy rule-labeled data into ChatGLM3 and
109 retrieve DBLP attributes via RAG.Zhong et al. (2024) run RoBERTa on cleaned citation contexts
110 and propagate node signals over a heterogeneous “abstract-title-reference” graph with GCN, then
111 ensemble.Chen et al. (2024b) prompt GPT-4/Claude in zero-shot with five prompt variants and
112 average their scores with LightGBM/CatBoost on structural features, no GPU required. However,
113 these top-performing approaches still carry inherent limitations. Grafting learning relies on a
114 cascade of separately trained models, amplifying error propagation and complicating hyper-parameter
115 tuning. The RAG pipeline further introduces retrieval noise that can dilute semantic focus. The
116 BERT-GCN hybrid demands meticulous graph construction and heavy feature engineering; its
117 performance drops when citation contexts are sparse or when the graph becomes overly dense. The
118 zero-shot LLM ensemble, despite its GPU-free advantage, still demands careful and labor-intensive
119 feature engineering—such as extracting citation frequencies, contextual keywords, and metadata—to
120 complement LLM outputs, limiting its scalability and adaptability across domains.

121 What’s more, the absence of standardized benchmarks has limited the comparability of these ap-
122 proaches. PST-Bench (Zhang et al., 2024) addresses this gap by providing a unified framework for
123 evaluating paper source tracing methods.

124 **3 PST Problem and PST-Bench Dataset**

125 **3.1 The Paper Source Tracing Problem**

126 Given a focal paper P and its full text, the goal is to identify its most salient references, here termed
127 *source papers*, which have significantly informed its ideas or methods. Each reference in P receives
128 an *importance weight* (0 to 1), quantifying its influence on the paper. The output for each paper P is
129 S_P .

130 A paper (p) may draw inspiration from one or more *primary sources*. A reference qualifies as a
131 *primary source* based on these criteria:

- 132 • Paper p ’s central idea is rooted in the reference.
133 • Paper p ’s fundamental methodology is derived from the reference.

134 **3.2 PST-Bench Dataset**

135 Given the inherent need for specialized domain knowledge in identifying paper sources, a cohort of
136 computer science graduate students was engaged to annotate papers within their respective areas of
137 expertise to build PST-Bench. The annotation workflow was structured as a collaborative online paper
138 group. Here, each student was tasked with presenting two papers weekly and identifying their corre-
139 sponding source papers. Following an extensive process involving data collection, rigorous expert



Figure 1: Muti_Agent Framework

140 validation, and thorough preprocessing, PST-Bench yielded a final collection of 2,141 meticulously
 141 labeled computer science publications.

142 4 Multi-Agent Ensemble for Paper Source Tracing

143 4.1 System Overview

144 We present a novel multi-agent ensemble architecture for paper source tracing, integrating Deepseek-
 145 R1-250528, GPT-5-2025-08-07, and Gemini-2.5-pro. The system operates through a structured six-
 146 component pipeline designed for robustness and scalability: XML Preprocessing, Prompt Engineering,
 147 Multi-Agent Prediction, Intelligent Ensemble, and Prediction Method. This architecture facilitates
 148 the precise identification of foundational methodological and conceptual papers within a citation
 149 network.

150 An illustrative diagram of the system's architecture is provided in Figure1. Next, we illustrate each
 151 component one by one.

152 4.2 XML Processing and Data Extraction

153 The system employs a robust dual-parsing strategy to accommodate diverse XML structures. A
 154 *primary parser*, leveraging `xml.etree`, handles standard TEI-compliant XML, ex-
 155 tracting reference titles, authors, publication venues, and years with high precision. A *fallback parser*,
 156 based on regular expressions, serves as a robust backup for malformed or non-standard XML. This is
 157 complemented by a comprehensive data cleaning pipeline that performs text normalization, including
 158 the removal of formatting characters, Unicode normalization, and whitespace standardization. This
 159 strategy ensures 100% reliability in processing a test set of 394 papers, with batch processing (100
 160 papers per batch) optimizing computational resource utilization.

161 4.3 Prompt Engineering Strategy

162 Our prompt engineering methodology is grounded in extensive empirical analysis, having evaluated
 163 over 1,000 human-annotated papers to identify the optimal prompt configuration achieving the highest
 164 F1-score. The unified prompt architecture incorporates multiple advanced reasoning frameworks
 165 derived from cutting-edge LLM research:

- **Counterfactual Reasoning:** Systematic evaluation of whether the target paper could be completed without each candidate reference
 - **Idea DNA Matching:** Identification of the earliest methodological and conceptual origins through citation chain analysis
 - **Multi-Role Socratic Dialogue:** Engagement of three distinct expert personas (Archaeologist, Experimentalist, Skeptic) in structured debate to reach consensus
- The prompt architecture enforces strict exclusion criteria to filter out non-source references including datasets, software tools, benchmarking studies, and general background literature, ensuring focused identification of true methodological source papers. All three models utilize customized optimized prompts and identical output formats to ensure consistency and comparability across predictions.

176 4.4 Multi-Agent Integration Strategy

177 We employ an advanced ensemble methodology that intelligently combines predictions from three
 178 complementary LLMs with optimized weight allocation. Empirically, we set the weight for each
 179 LLM as follows: DeepSeek-R1-0528 (0.3), GPT-5 (0.35), Gemini-2.5-pro (0.35).

180 4.4.1 Confidence Score Extraction

181 The system implements a multi-format confidence extraction mechanism supporting both structured
 182 JSON output and legacy formats:

```
183 {"source_references": [], "confidence_scores": {}, "reasoning": "detailed analysis"}
```

184 Note that sometimes the LLMs do not output confidence scores.

185 4.4.2 Intelligent Default Scoring

186 When explicit confidence scores are unavailable, the scoring method assigns a base score of 0.3 to
 187 each mentioned reference, augmenting it by 0.2 for inclusion in the prediction from DeepSeek-R1,
 188 GPT, or Gemini individually, plus an additional 0.1 bonus if present in all three, with the final score
 189 capped at 1.0.

```
190 1 scores = []
191 2 for bid in all_sources:
192 3     base_score = 0.3
193 4     is_in_ds = bid in ds_sources
194 5     is_in_gpt = bid in gpt_sources
195 6     is_in_gemini = bid in gemini_sources
196 7
197 8     # Add points for each source type if present
198 9     base_score += (is_in_ds * 0.2)
199 10    base_score += (is_in_gpt * 0.2)
200 11    base_score += (is_in_gemini * 0.2)
201 12
202 13    # Add bonus if in all three
203 14    base_score += (is_in_ds and is_in_gpt and is_in_gemini) * 0.1
204 15
205 16    scores[bid] = min(1.0, base_score)
```

Listing 1: Python Implementation of Ensemble Scoring

206 4.4.3 Consistency Penalty Mechanism

207 A dynamic penalty function is applied based on maximum pairwise score differences between models,
 208 with penalty factors ranging from 0.1 (maximal disagreement) to 1.0 (minimal difference) to mitigate
 209 inconsistent predictions.

210 **4.4.4 Probability Distribution Conversion**

211 The penalty score for the i -th LLM is calculated as follows:

$$P(i) = C(|s_{deepseek} - s_{gpt}|) + C(|s_{gpt} - s_{gemini}|) + C(|s_{gemini} - s_{deepseek}|)$$

212 where D is the Consistency Penalty Factor and s_{model} is the prediction probability of the current
213 LLM.

214 **4.5 Prediction Method**

215 Finally, the prediction score of the i -th LLM is defined as:

$$y_i = P(i) * w(i) * s(i)$$

216 where $P(i)$ is the penalty score, $w(i)$ is the weight score of the LLM, and $s(i)$ is the initial prediction
217 score of the LLM.

218 **5 Evaluation Framework**

219 **5.1 Metrics**

220 We adopt mean average precision (MAP) to evaluate the prediction results. Concretely, for each
221 paper p in the test set,

$$\text{AP}(p) = \frac{1}{R_p} \sum_{k=1}^{M_p} \text{Prec}_p(k) 1_k, \quad (1)$$

222 where R_p is the number of reference sources of paper p , M_p is the number of references of paper p ,
223 $\text{Prec}_p(k)$ is the precision at cut-off k in the ranked output list $S_p(k)$, and 1_k is the indicator function
224 for the k -th item being a relevant document, with the values 0 or 1.

$$\text{MAP} = \frac{1}{|\mathcal{P}_{\text{test}}|} \sum_{p \in \mathcal{P}_{\text{test}}} \text{AP}(p), \quad (2)$$

225 where $\mathcal{P}_{\text{test}}$ is the set of papers in the testing set.

226 **5.2 Baseline Models**

227 We implement several baseline models for comparison to validate the design of our approach:
228 DeepSeek-R1-0528, GPT-5-2025-08-07, and Gemini-2.5-pro.

229 **6 Experiments and Results**

230 **6.1 Experimental Setup**

231 The system is implemented in Python 3.8+. API configurations are centrally managed through a
232 structured config.py file, supporting all three model providers with proper authentication, endpoint
233 configuration, and rate limit management. The implementation includes comprehensive unit testing,
234 integration testing, and end-to-end validation procedures. The system architecture follows a modular
235 design with separate components for XML processing, model integration, quality assessment, and
236 submission generation, enabling maintainability and future extensibility.

237 For evaluation, we utilize the PST-Bench dataset (Zhang et al., 2024), which comprises 1,576
238 meticulously labeled computer science papers. The dataset is partitioned based on publication year,
239 with 788 papers allocated for training, 394 for validation, and the remaining 394 reserved for testing.
240 This temporal split ensures that models are evaluated on papers published after those in the training
241 set, simulating real-world deployment scenarios.

Final Leaderboard - PST-KDD-2024

If you see anyone use multiple accounts, please let us know.

#	Team Name	Score	#Parameters	Total GPU Memory (Gb)	Final Entries
1	jkassieslzh	0.51318	8,000,000,000	640	153
2	BlackPearl	0.48789	6,000,000,000	640	33
3	NJUST_KMG	0.47961	125,000,000	24	34
4	pst-auto-agent	0.46936	0	0	5
5	英国大力士	0.46121	0	0	21
6	chinesegept	0.44492	440,000,000	24	34
7	listenn	0.44492	0	0	1
8	😊data	0.44278	560,000,000	80	20

Figure 2: Leaderboard of KDD Cup 2024.

242 6.2 Main Results

243 Table 1 shows the performance of baseline models on PST-Bench:

Table 1: Performance of baseline models on PST-Bench

Model	MAP
DeepSeek-R1-250528	0.246
GPT-5-2025-08-07	0.315
Gemini-2.5-pro	0.318
pst-auto-agent	0.388

244 6.3 Analysis

245 The DeepSeek-R1-250528 model exhibited the lowest performance, achieving a MAP score of 0.246.
246 This positions it as the weakest baseline within this comparison.

247 The GPT-5-2025-08-07 and Gemini-2.5-pro models demonstrated significantly improved performance
248 relative to DeepSeek-R1-250528, with MAP scores of 0.315 and 0.318, respectively. The marginal
249 difference between these two models (0.003) suggests comparable efficacy on the PST-Bench dataset.

250 The `pst-auto-agent` model achieved the highest performance, with a MAP score of 0.388. This
251 represents a notable advancement, surpassing Gemini-2.5-pro by approximately 22.0% relatively.
252 The superior performance of `pst-auto-agent` suggests its architectural design confers a substantial
253 advantage for the PST-Bench task.

254 6.4 Further Results on KDD Cup 2024

255 The proposed method is a tuning-free method that do not require feature engineering. We further
256 enhance the top-ranked method **English Hercules** in KDD Cup 2024¹ by ensembling our method into
257 its framework. Generally speaking, English Hercules is a GPU-free approach the combines feature
258 engineering and LLM API-based methods. To this end, we integrate GPT-5, DeepSeek-R1, and
259 Gemini-pro into its framework by utilizing the ensembling method of English Hercules. Our method
260 achieved 4th place overall on the KDD Cup 2024 leaderboard and ranked 1st among all GPU-free
261 methods. The results are shown in Figure 2, demonstrating that our approach clearly complements
262 their feature engineering and LLM-based approaches.

¹https://www.biendata.xyz/competition/pst_kdd_2024/final-leaderboard/

263 **7 Conclusion and Future Work**

264 This paper proposed a novel multi-agent ensemble architecture, termed `pst-auto-agent` for the
265 paper source tracing task, which integrates Deepseek-R1-250528, GPT-5-2025-08-07, and Gemini-
266 2.5-pro within a structured pipeline. This architecture incorporates advanced XML preprocessing,
267 empirically optimized prompt engineering, a sophisticated multi-agent prediction mechanism, and an
268 intelligent ensemble strategy that includes confidence scoring and a consistency penalty.

269 Experimental results on PST-Bench demonstrated that `pst-auto-agent` achieved a Mean Average
270 Precision (MAP) score of 0.388, significantly outperforming individual baseline models. Furthermore,
271 when integrated with the top-ranked English Hercules framework in KDD Cup 2024, our method
272 exhibited a complementary effect, enhancing overall performance. This work underscores the efficacy
273 of a multi-agent ensemble approach for the challenging task of identifying primary source papers.

274 Future work could explore several directions:

- 275 • Incorporating temporal dynamics of citation patterns
276 • Developing domain-specific adaptations for different research areas
277 • Exploring interactive tools for researchers to explore citation networks

278 **References**

279 David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine*
280 *Learning research*, 3(Jan):993–1022, 2003.

281 Haoru Chen, Xiaocheng Zhang, Yang Zhou, Mengjiao Bao, and Peng Yan. Grafting learning is all
282 you need. In *KDD 2024 Workshop OAGChallenge Cup*, 2024a. URL <https://openreview.net/forum?id=s2hVfSPTbU>.

284 Kunlong Chen, Junjun Wang, Zhaoqun Chen, Kunjin Chen, and Yitian Chen. Llm-powered ensemble
285 learning for paper source tracing: A gpu-free approach. In *KDD 2024 Workshop OAG-Challenge*
286 *Cup*, 2024b. URL <https://openreview.net/forum?id=H5wbKuYjxR>.

287 Daryl E. Chubin and Eugene Garfield. Is citation analysis a legitimate evaluation tool? *Scientometrics*,
288 2(1):91–94, 1980. doi: 10.1007/BF02016602. URL <https://doi.org/10.1007/BF02016602>.

289 MM KESSLER. Bibliographic coupling between scientific papers. 14:10–25, 1963.

290 Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks.
291 In *International Conference on Learning Representations*, 2017.

292 Quoc V. Le and Tomas Mikolov. Distributed representations of sentences and documents.
293 [abs/1405.4053](https://arxiv.org/abs/1405.4053), 2014.

294 Henry Small. Co-citation in the scientific literature: A new measure of the relationship between
295 two documents. *J. Am. Soc. Inf. Sci.*, 24(4):265–269, 1973. doi: 10.1002/ASI.4630240406. URL
296 <https://doi.org/10.1002/asi.4630240406>.

297 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
298 Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information*
299 *processing systems*, pp. 5998–6008, 2017.

300 Fanjin Zhang, Kun Cao, Yukuo Cen, Jifan Yu, Da Yin, and Jie Tang. Pst-bench: Tracing and
301 benchmarking the source of publications. In *Computing Research Repository*, 2024.

302 Shupeng Zhong, Xinger Li, Shushan Jin, and Yang Yang. The solution for the `pst-kdd-2024 oag-`
303 `challenge`. In *KDD 2024 Workshop OAGChallenge Cup*, 2024. URL <https://openreview.net/forum?id=0v0EZ7PuuI>.

305 **Agents4Science AI Involvement Checklist**

- 306 1. **Hypothesis development:** Hypothesis development includes the process by which you
307 came to explore this research topic and research question. This can involve the background
308 research performed by either researchers or by AI. This can also involve whether the idea
309 was proposed by researchers or by AI.

310 Answer: **[A]**

311 Explanation: The research hypothesis and questions were developed entirely by human
312 researchers based on existing literature and identified research gaps in paper source tracing.

- 313 2. **Experimental design and implementation:** This category includes design of experiments
314 that are used to test the hypotheses, coding and implementation of computational methods,
315 and the execution of these experiments.

316 Answer: **[C]**

317 Explanation: AI assistance designed the experimental framework and implemented the
318 core algorithms coding, with Human researchers participated in code optimization and data
319 processing tasks.

- 320 3. **Analysis of data and interpretation of results:** This category encompasses any process to
321 organize and process data for the experiments in the paper. It also includes interpretations of
322 the results of the study.

323 Answer: **[B]**

324 Explanation: Data analysis was primarily conducted by human researchers, with AI tools
325 used for statistical calculations and visualization. Interpretation of results was entirely
326 human-driven.

- 327 4. **Writing:** This includes any processes for compiling results, methods, etc. into the final
328 paper form. This can involve not only writing of the main text but also figure-making,
329 improving layout of the manuscript, and formulation of narrative.

330 Answer: **[C]**

331 Explanation: The writing process involved significant AI assistance for drafting and editing,
332 while human researchers provided overall structure, critical analysis, and final review.

- 333 5. **Observed AI Limitations:** What limitations have you found when using AI as a partner or
334 lead author?

335 Description: AI systems demonstrated limitations in understanding long complex academic
336 context, maintaining consistent technical accuracy, and providing novel insights beyond
337 pattern recognition from training data.

338 **Agents4Science Paper Checklist**

339 **1. Claims**

340 Question: Do the main claims made in the abstract and introduction accurately reflect the
341 paper's contributions and scope?

342 Answer: [Yes]

343 Justification: The abstract and introduction clearly state our contributions including the
344 multi-agent ensemble architecture, experimental results, and integration with KDD Cup
345 2024 top solution.

346 Guidelines:

- 347 • The answer NA means that the abstract and introduction do not include the claims
348 made in the paper.
- 349 • The abstract and/or introduction should clearly state the claims made, including the
350 contributions made in the paper and important assumptions and limitations. A No or
351 NA answer to this question will not be perceived well by the reviewers.
- 352 • The claims made should match theoretical and experimental results, and reflect how
353 much the results can be expected to generalize to other settings.
- 354 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
355 are not attained by the paper.

356 **2. Limitations**

357 Question: Does the paper discuss the limitations of the work performed by the authors?

358 Answer: [Yes]

359 Justification: Section 6 discusses limitations and future work directions including domain-
360 specific adaptations and computational efficiency considerations.

361 Guidelines:

- 362 • The answer NA means that the paper has no limitation while the answer No means that
363 the paper has limitations, but those are not discussed in the paper.
- 364 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 365 • The paper should point out any strong assumptions and how robust the results are to
366 violations of these assumptions (e.g., independence assumptions, noiseless settings,
367 model well-specification, asymptotic approximations only holding locally). The authors
368 should reflect on how these assumptions might be violated in practice and what the
369 implications would be.
- 370 • The authors should reflect on the scope of the claims made, e.g., if the approach was
371 only tested on a few datasets or with a few runs. In general, empirical results often
372 depend on implicit assumptions, which should be articulated.
- 373 • The authors should reflect on the factors that influence the performance of the approach.
374 For example, a facial recognition algorithm may perform poorly when image resolution
375 is low or images are taken in low lighting.
- 376 • The authors should discuss the computational efficiency of the proposed algorithms
377 and how they scale with dataset size.
- 378 • If applicable, the authors should discuss possible limitations of their approach to
379 address problems of privacy and fairness.
- 380 • While the authors might fear that complete honesty about limitations might be used by
381 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
382 limitations that aren't acknowledged in the paper. Reviewers will be specifically
383 instructed to not penalize honesty concerning limitations.

384 **3. Theory assumptions and proofs**

385 Question: For each theoretical result, does the paper provide the full set of assumptions and
386 a complete (and correct) proof?

387 Answer: [NA]

388 Justification: This paper focuses on empirical benchmarking and experimental results, with
389 no theoretical results requiring formal proofs.

390 Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.

398 **4. Experimental result reproducibility**

399 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
400 perimental results of the paper to the extent that it affects the main claims and/or conclusions
401 of the paper (regardless of whether the code and data are provided or not)?

402 Answer: [Yes]

403 Justification: Section 5 describes the experimental setup, evaluation metrics, and implemen-
404 tation details necessary for reproducibility.

405 Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

416 **5. Open access to data and code**

417 Question: Does the paper provide open access to the data and code, with sufficient instruc-
418 tions to faithfully reproduce the main experimental results, as described in supplemental
419 material?

420 Answer: [Yes]

421 Justification: The PST-Bench dataset and code implementation details are described in the
422 paper, with access information provided in the documentation.

423 Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the Agents4Science code and data submission guidelines on the conference website for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).

435 **6. Experimental setting/details**

436 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
437 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
438 results?

439 Answer: [Yes]

440 Justification: Section 5 describes the experimental setup, including dataset details, evaluation
441 metrics (MAP), and implementation specifics with Python 3.8+ and API configurations.

442 Guidelines:

- 443 • The answer NA means that the paper does not include experiments.
- 444 • The experimental setting should be presented in the core of the paper to a level of detail
- 445 that is necessary to appreciate the results and make sense of them.
- 446 • The full details can be provided either with the code, in appendix, or as supplemental
- 447 material.

448 **7. Experiment statistical significance**

449 Question: Does the paper report error bars suitably and correctly defined or other appropriate
450 information about the statistical significance of the experiments?

451 Answer: [No]

452 Justification: The paper reports MAP scores but does not include error bars or statistical
453 significance tests due to the nature of the single evaluation run on the test set.

454 Guidelines:

- 455 • The answer NA means that the paper does not include experiments.
- 456 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
- 457 dence intervals, or statistical significance tests, at least for the experiments that support
- 458 the main claims of the paper.
- 459 • The factors of variability that the error bars are capturing should be clearly stated
- 460 (for example, train/test split, initialization, or overall run with given experimental
- 461 conditions).

462 **8. Experiments compute resources**

463 Question: For each experiment, does the paper provide sufficient information on the com-
464 puter resources (type of compute workers, memory, time of execution) needed to reproduce
465 the experiments?

466 Answer: [No]

467 Justification: The paper mentions using commercial LLM APIs (DeepSeek, GPT-5, Gemini)
468 but does not provide detailed compute resource information for local processing components.

469 Guidelines:

- 470 • The answer NA means that the paper does not include experiments.
- 471 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
- 472 or cloud provider, including relevant memory and storage.
- 473 • The paper should provide the amount of compute required for each of the individual
- 474 experimental runs as well as estimate the total compute.

475 **9. Code of ethics**

476 Question: Does the research conducted in the paper conform, in every respect, with the
477 Agents4Science Code of Ethics (see conference website)?

478 Answer: [Yes]

479 Justification: The research focuses on academic paper source tracing for knowledge discov-
480 ery and follows ethical guidelines for AI research and citation analysis.

481 Guidelines:

- 482 • The answer NA means that the authors have not reviewed the Agents4Science Code of
- 483 Ethics.
- 484 • If the authors answer No, they should explain the special circumstances that require a
- 485 deviation from the Code of Ethics.

486 **10. Broader impacts**

487 Question: Does the paper discuss both potential positive societal impacts and negative
488 societal impacts of the work performed?

489 Answer: [No]

490 Justification: The paper focuses on technical contributions to paper source tracing but does
491 not explicitly discuss broader societal impacts in a dedicated section.

492

Guidelines:

493

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations, privacy considerations, and security considerations.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies.

494

495

496

497

498

499

500