

Paper Circle: An Open-source Multi-agent Research Discovery and Analysis Framework

Anonymous ACL submission

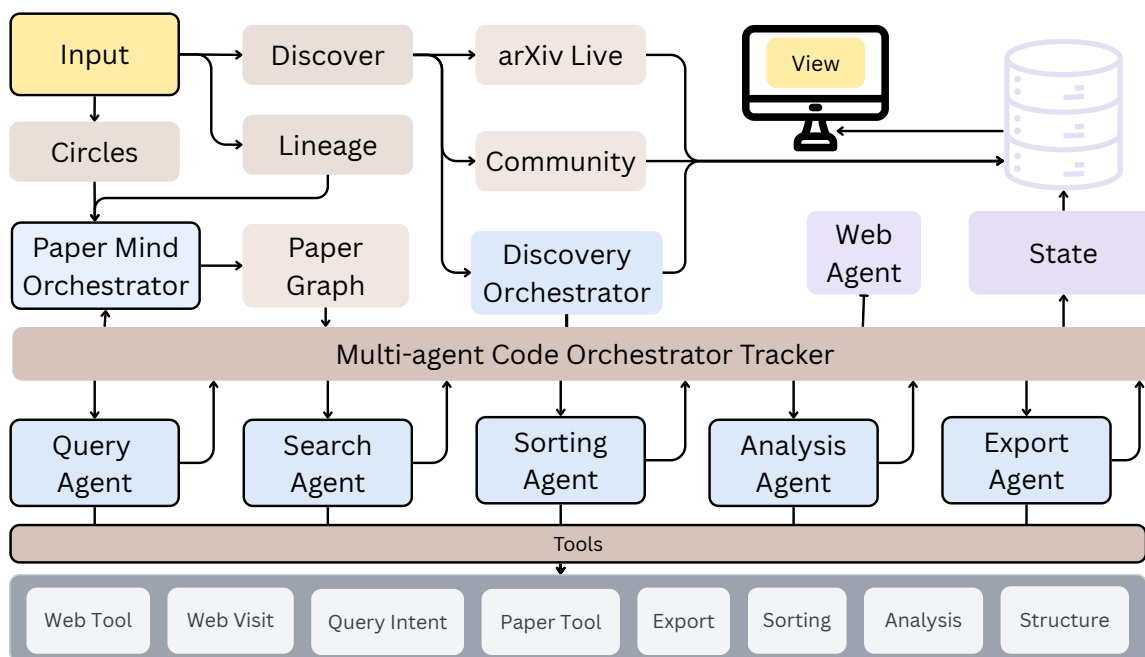


Figure 1: Overview of the Paper Circle pipeline. Given a user query, Paper Circle builds a paper set from multiple sources (e.g., paper graph, community, and arXiv live) via the Paper Mind for analysis and Discovery Orchestrators for search of the paper. A multi-agent layer (query, search, sorting, analysis, export) is coordinated by the Tracker, which maintains a shared state that is persisted to a backing store and displayed to the user through interface.

Abstract

The rapid growth of scientific literature has made it increasingly difficult for researchers to efficiently discover, evaluate, and synthesize relevant work. Recent advances in multi-agent large language models (LLMs) have demonstrated strong potential for understanding user intent and are being trained to utilize various tools. In this paper, we introduce Paper Circle, a multi-agent research discovery and analysis system designed to reduce the effort required to find, assess, organize, and understand academic literature. The system comprises two complementary pipelines: (1) a Discovery Pipeline that integrates offline and online retrieval from multiple sources, multi-criteria scoring, diversity-aware ranking, and structured outputs; and (2) an Analysis Pipeline that

transforms individual papers into structured knowledge graphs with typed nodes (e.g., concepts, methods, experiments, and figures) and edges, enabling graph-aware question answering and coverage verification. Both pipelines are implemented within a coder LLM-based multi-agent orchestration framework and produce fully reproducible, synchronized outputs (JSON, CSV, BibTeX, Markdown, and HTML) at each agent step. This paper describes the system architecture, agent roles, retrieval and scoring methods, knowledge graph schema, and evaluation interfaces that together form the Paper Circle research workflow. We benchmark Paper Circle on both paper retrieval and paper review generation, reporting hit rate, MRR, and Recall@K. Results show consistent improvements with stronger agent models. We will

publicly release the website and source code.

lighting strengths and weaknesses to guide human reading priorities (Naumov et al., 2025).

1 Introduction

The pace of scientific publication has accelerated exponentially, creating a significant burden on researchers attempting to stay abreast of new developments (Reddy and Shojaee, 2025; Pramanick et al., 2023). Traditional search engines and recommendation systems often struggle to provide the depth and context required for rigorous literature reviews, leading to fragmented discovery workflows. Recently, the advent of Large Language Models (LLMs) has catalyzed a shift towards "AI Scientists", autonomous multi-agent systems (MAS) capable of generating hypotheses, conducting experiments, and even writing papers (Chen et al., 2025b; Naumov et al., 2025). While these systems demonstrate the potential of agentic workflows, there remains a critical gap between fully autonomous simulations and the practical, collaborative needs of human research communities.

Paper Circle addresses (as shown in the Figure 1) this gap by introducing a comprehensive *Multi-Agent Research Platform* that supports the entire lifecycle of literature engagement: from discovery and analysis to critique and synthesis. Unlike purely autonomous systems that aim to replace the researcher, Paper Circle is designed as a collaborative workbench that augments human intelligence through three integrated subsystems:

- 1. Discovery Pipeline:** A multi-agent retrieval system that goes beyond simple keyword matching. It employs a multi-dimensional scoring framework to surface high-value research. Crucially, this pipeline is deterministic and produces structured artifacts (JSON, linear logs) at every step.
- 2. Paper Mind Graph:** To facilitate deep understanding, Paper Circle constructs a dynamic Knowledge Graph from retrieved literature. This "Paper Mind" enables researchers to query the collective intelligence of a reading list, identifying latent connections between disparate works and supporting complex Question-Answering workflows that are grounded in specific citation sub-graphs.
- 3. Review Agents:** This platform features a team of specialized review agents that generate detailed critiques and scores, consistently high-

By integrating these capabilities into a shared "Reading Circle" environment, Paper Circle transforms literature review from a solitary task into a community-driven, AI-augmented operation.

2 Related Work

2.1 Autonomous Scientific Discovery

The emerging field of AI-Scientists aim to automate the entire research lifecycle. Systems like DORA AI agent (Naumov et al., 2025) and EvoResearch (Gajjar, 2025) demonstrate end-to-end capabilities, from hypothesis generation to report writing. Similarly, O-Research (Li et al., 2025), MARS (Chen et al., 2025a), and AlphaResearch (Yu et al., 2025c) treat research as a multi-step optimization problem, often using reinforcement learning to refine discovery strategies. Specialized agents have also been proposed for causal discovery, such as CausalSteward (Wang et al., 2025) and other multi-agent frameworks (Le et al., 2025). While these systems push the boundaries of autonomy, Paper Circle prioritizes *curation and reproducibility* over full automation. Instead of replacing the researcher, Paper Circle acts as a force multiplier for human teams, ensuring that the discovery process remains transparent and verifiable.

2.2 MAS in Specialized Domains

MAS have shown remarkable success in specific scientific verticals. In chemistry and materials science, frameworks like ChemThinker (Ju et al., 2025), MOOSE-Chem (Yang et al., 2025), and ChemBOMAS (Han et al., 2025a) leverage LLMs to discover new molecules and optimize experiments (Kumbhar et al., 2025). In biology and healthcare, agents facilitate single-cell analysis (CellAgent (Xiao et al., 2024)), phenotype discovery (PhenoGraph (Niyakan and Qian, 2025)), and clinical data analysis (Spieser et al., 2025). Other applications range from drug discovery (Fehlis et al., 2025) and psychiatry diagnosis (Xiao et al., 2025) to financial forecasting, where systems like ASTRAFIN (Singh and Kumar, 2025) and other stock analysis agents (Chandrashekar et al., 2025; Wawer and Chudziak, 2025) predict market trends. Paper Circle complements these domain-specific tools by providing a *general-purpose* discovery pipeline that can be adapted to any discipline, serv-

ing as the foundational layer for literature review and knowledge management.

2.3 Community Simulation and Collaboration

A distinct line of research focuses on simulating or facilitating the social aspects of science. Research-Town (Yu et al., 2025a,b) models the research community using agents to understand how ideas propagate. Other works explore collaborative dynamics through automated negotiation (NegoLog (Doğru et al., 2024), NEGOTIATOR (Keskin et al., 2024)) and cohesive dialogue generation (Chu et al., 2024). Frameworks like PiFlow (Pu et al., 2025), RED-EREF (Yuan and Xie, 2025), and blackboard systems (Salemi et al., 2025) propose mechanisms for agent collaboration in information discovery. Paper Circle distinguishes itself by moving beyond simulation; it provides a real-world platform for *human-AI collaboration*. It does not just model how researchers might interact, but actively facilitates those interactions through shared reading lists, discussion threads, and collaborative ranking.

3 Methodology

3.1 Background

Multi-Agent Systems (MAS) represent a paradigm where autonomous entities interact to solve complex problems distributedly. In the context of scientific discovery, MAS allows for the decomposition of intricate research tasks, such as literature search, reading, and reasoning, into manageable sub-routines handled by specialized agents (Wooldridge, 2002). Unlike monolithic LLM approaches, agentic workflows can maintain distinct personas (e.g., "The Skeptic", "The Creative") and leverage external tools, reducing hallucination and improving reasoning depth through inter-agent dialogue (Reddy and Shojaee, 2025).

The baseline for our orchestration layer is the smolagents (Roucher et al., 2025) library. The pipeline uses a CodeAgent (CoA) as the central orchestrator, which can attend parallel agent calls and toll calls and multiple ToolCallingAgent (ToCA) instances, each attached to specific capabilities (e.g., arXiv retrieval, PDF parsing). The baseline responsibilities include (i) tool invocation, (ii) multi-step planning via the orchestrator, and (iii) delegation to specialized agents. PaperCircle extends this foundation by adding structured outputs, offline search capabilities, and rigorous evaluation metrics. We preserve the baseline tool interface,

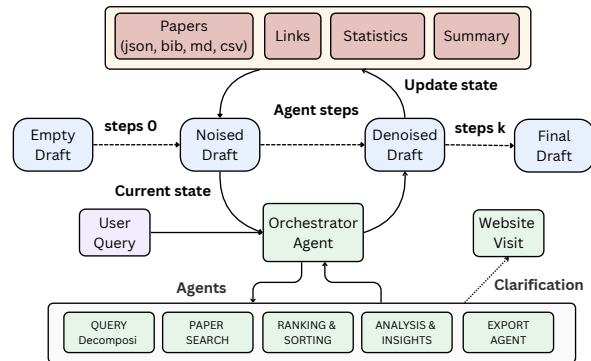


Figure 2: The main iterative diagram for the paper discovery framework. The system maintains an explicit, evolving discovery state (papers, links, statistics, and summaries) that is iteratively updated through agentic steps. Starting from an empty draft, the orchestrator agent alternates between noising and denoising operations over multiple steps, progressively refining the draft into a final result. When necessary, a web search agent is invoked for clarification or recent information.

where each tool receives explicit parameters and returns a formatted string response, allowing the orchestrator to chain steps while maintaining high readability and traceability.

3.2 System Architecture

Figure 1 illustrates the overall architecture of Paper Circle. The system consists of two complementary multi-agent pipelines: the *Discovery Pipeline* for finding relevant papers, and the *Analysis Pipeline* for deep understanding of individual papers.

3.3 Paper Discovery Agent Design

The main diagram of the discovery subsystem is shown in Figure 2, which is composed of multiple agents, each bound to a small, explicit tool interface. It is inspired by the TTD-DR (Han et al., 2025b) for iteratively updating the updated version at each agentic step. The core agents are:

Intent Classification Agent. Parses user text into search mode (offline, online, or both), conference filters, year range, and ranking preferences. Most importantly, it uses a web agent in the pipeline for any unclear queries or recent knowledge.

Paper Search Agent. Executes offline or online retrieval based on intent, merges results, performs deduplication, and updates state and outputs.

Sorting Agent. Reorders papers using recency, citations, similarity, novelty, BM25 (Chen and Wiseman, 2023), or combined weights; or applies a cross-encoder reranker (Wang et al., 2020).

Analysis Agent. Computes aggregate statistics

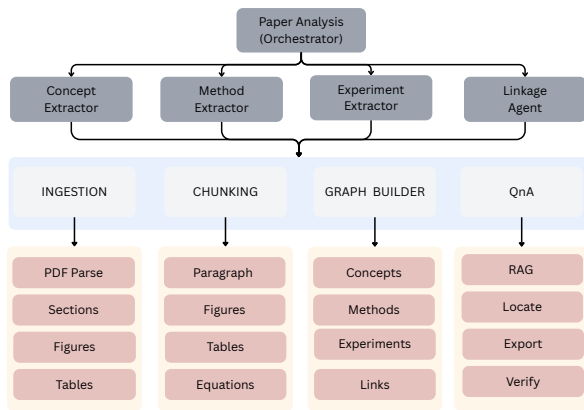


Figure 3: A paper analysis orchestrator agents for concepts, methods, experiments, and cross-entity linkages. The pipeline consists of four main stages: ingestion, which parses PDFs into structured elements (sections, figures, tables, equations); semantic chunking, which produces structure-aware text units; graph construction, which builds a typed knowledge graph of concepts, methods, experiments, and their relations with full traceability to source text; and a Q&A layer that enables graph-aware retrieval, verification, and export.

and insights, including source distribution, year trends, and top authors.

Export Agent. Produces synchronized exports and provides a consistent interface for downstreaming.

Web Search Agent. Provides auxiliary access to web search tools when online lookups are required.

3.4 Paper Analysis Agent

While the discovery pipeline addresses the challenge of finding relevant papers, researchers also need to understand and synthesize the content of individual papers deeply (Korat, 2025). Paper Circle addresses this with a complementary *Paper Analysis Agent* that transforms research papers into structured, queryable knowledge graphs with full traceability to the original text. The Paper Analysis Agent operates as a multi-stage pipeline with four specialized components as shown in the figure: (1) Ingestion Layer, (2) Graph Builder, (3) Q&A System, and Verification Layer.

PDF Ingestion and Chunking. The ingestion pipeline uses PyMuPDF for robust PDF parsing (Adhikari and Agarwal, 2024). The PDFParser class extracts: **Metadata:** Title, authors, abstract, arXiv ID, venue, and page count. **Sections:** Hierarchical section structure with parent-child relationships, identified via numbering patterns (e.g., “1.2 Background”). **Figures and Tables:** Caption

text, page locations, and nearby context for linkage. **Equations:** Numbered equations with surrounding context.

Unlike token-based chunking, the SemanticChunker (Qu et al., 2025) creates chunks aligned with document structure. Paragraphs within sections are grouped up to a configurable limit (default 1500 characters), while figures, tables, and equations are preserved as distinct chunks with their captions and context.

Knowledge Graph Schema. The mind graph follows a typed schema with nodes (Zhang et al., 2025a) for papers, sections, concepts, methods, experiments, datasets, and visual elements (figures, tables, equations), and edges encoding structural and semantic relations (e.g., hierarchy, definition, proposal, usage, evaluation, illustration, dependency). All nodes and edges carry provenance metadata—including source chunk IDs, page numbers, verification status, confidence scores, and timestamps—ensuring full traceability to the original PDF.

3.5 Multi-Agent Extraction

The GraphBuilder (Zhu et al., 2024b) orchestrates four specialized CoA-based extractors. The *Concept Extractor* identifies and classifies key concepts by type and importance; the *Method Extractor* extracts algorithms and techniques from method sections; the *Experiment Extractor* recovers experimental setups, datasets, metrics, and results; and the *Linkage Agent* connects figures and tables to the concepts or methods they illustrate. Extraction proceeds in staged phases—concepts, methods, experiments, visual linkage, and inter-concept relations—each incrementally updating the shared MindGraph.

Graph-Aware Q&A. The Q&A module combines vector retrieval with graph traversal. An EmbeddingStore indexes text chunks and node descriptions, while the GraphRetriever retrieves top- k relevant nodes and chunks and expands context via 1-hop neighbors. The PaperQA agent generates answers grounded in retrieved text, graph relations, and linked figures or tables, and returns supporting evidence with confidence estimates. A locate function enables precise localization of concepts, figures, or tables by page and context.

Coverage Verification. To prevent silent omissions, a CoverageChecker evaluates figure, table, section, and equation coverage, producing an over-

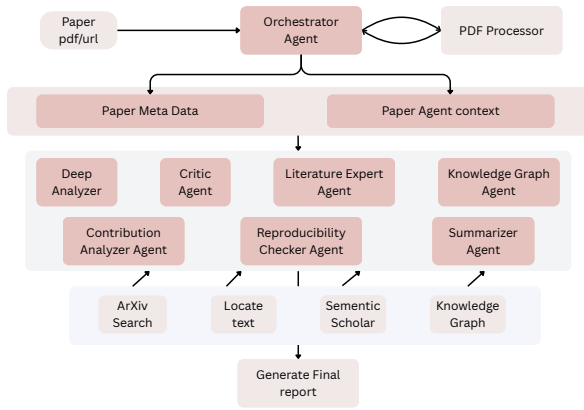


Figure 4: Multi-agent paper analysis and review architecture. Given a paper specified by a PDF or URL, an orchestrator agent coordinates PDF processing and maintains shared paper metadata and agent context. Specialized agents operate in parallel to perform deep technical analysis, contribution extraction, critical review, literature linking, reproducibility checking, summarization, and knowledge graph construction. External tools such as arXiv search, Semantic Scholar, and targeted text localization are invoked as needed. The orchestrator aggregates agent outputs into a unified, structured final report, enabling comprehensive, reviewer-style analysis with modular extensibility.

all coverage score and identifying unlinked or missing elements with actionable diagnostics. This provides a lightweight quality assurance step prior to downstream use.

3.6 Research Review Framework

In Sec. 3.4, we describe the paper analysis of agentic capabilities, which we further extend for automated peer-review-style assessment. Unlike AgentReview (Jin et al., 2024; D’Arcy et al., 2024), we follow the paper analysis perspective, which not only provides the review but also builds a strong graph between the concepts.

Architecture. The system is built upon a multi-agent orchestration framework (Figure 4) that coordinates the execution of seven specialized roles. Each agent is instantiated as a ToCA or CoA (Roucher et al., 2025).

Deep Analyzer. Focuses on the technical core of the paper. It breaks down the mathematical foundations, identifies specific methodology components, and extracts primary experimental findings.

Critic. Emulates a senior conference reviewer (e.g., NeurIPS, ICML). It provides a rigorous assessment of strengths and weaknesses, generates

author-facing questions, and assigns scores for novelty, clarity, and significance.

Literature Expert. Interfaces with external academic databases including Semantic Scholar and arXiv. It maps the paper’s position within the existing research landscape and verifies citation accuracy.

Contribution Analyzer. Separates explicit author claims from verified technical contributions, identifying potential overclaiming or missing baseline comparisons.

Reproducibility Checker: Quantifies the transparency of the research by assessing the availability of source code, hyperparameter specifications, dataset accessibility, and compute requirement disclosures.

Summarizer. Generates multi-fidelity summaries across different abstraction levels, ranging from concise executive summaries to deep technical precis.

Orchestration and Pipeline Execution The Multi Agent Orchestrator manages the lifecycle of these agents through a multi-stage pipeline. The system supports parallel execution using a ThreadPoolExecutor.

4 Experiments

4.1 Experimental setup

All the experiments are done with open-source model with 4×40 GB Nvidia GPUs. We used the Ollama¹ platform with the fastllm library (Gong et al., 2025).

Database Curation. We curated a diverse corpus, as shown in Table 1 of research papers from leading CS and ML conferences, primarily sourced from OpenReview² and augmented with metadata and peer-review information.

Evaluation. Paper Circle provides built-in evaluation metrics. When a ground-truth paper title or identifier is provided, the system computes Mean Reciprocal Rank (MRR), Recall@K, Precision@K, and hit rates. These metrics are computed per step and stored in the JSON file for longitudinal tracking. For batch evaluation, a parallel benchmarking utility executes multiple queries concurrently and aggregates mean metrics and timing statistics. This

¹<https://ollama.com/>

²<https://openreview.net/>

Conference	ICLR	NeurIPS	ICML	CVPR	IROS	ICRA	AAAI	ACL	ICCV	EMNLP	Other
Count	12	39	13	13	25	25	5	5	7	4	144

Table 1: The Database corpus across major conferences. The “Other” category includes venues such as AISTATS, RSS, SIGGRAPH, and WACV. **Count** indicates the number of the most recent conference venue included.

supports lightweight comparisons between search configurations (offline vs. online, BM25 (Chen and Wiseman, 2023) vs. semantic (all-MiniLM-L6-v2 (Wang et al., 2020)), with or without Qwen3-Reranker-0.6B (Zhang et al., 2025b)) without requiring external tooling.

Baseline Agent. This framework is developed using the Smolagent multi-agent tool, calling the (ToCA) agent and the code agent (CoA), with tools utilized being manually developed.

Architecture. We evaluate multiple retrieval baselines: bm25, bm25+reranker (BM25 (Chen and Wiseman, 2023)& cross-encoder (Zhang et al., 2025b)), reranker (Zhang et al., 2025b), semantic (Wang et al., 2020), and hybrid (BM25 combined with semantic retrieval). We also compare pipeline structures with different agent compositions: full includes all five agents (intent, search, sort, analysis, export), minimal uses only the search agent, search_sort uses search and sort, search_analysis uses search and analysis, and no_intent is a full pipeline with no intent.

4.2 Results

Natural Text-based retrieval. We evaluate our multi-agent paper retrieval system across multiple LLMs and retrieval baselines. We did two query type experiments, one a research assistant-based natural queries generated by running gpt-oss-20B models (called RABench), and randomly sampling one paper record from the database, extracting a concise “topic” phrase from its title, keywords or abstract, then picking a natural-language template and optional prefix to turn that topic into a realistic search query. We also randomly chose a scope (conference/year/range/none) to add corresponding text to the query and to emit matching structured filters. This query we referred to as SemanticBench.

All experiments were conducted on a 50 query benchmark, measuring the success rate, the hit rate, the mean reciprocal rank (MRR), and the recall.

Model Comparison. Table 2 presents comprehensive evaluation results comparing agent-based

models with retrieval baselines. The results reveal a clear performance hierarchy across methods and scales. Two agent models achieve the highest retrieval effectiveness with an 80% hit rate, qwen3-coder-30b-Q3KM (quantized) and qwen3-coder:30b—with qwen3-coder-30b-Q3KM also delivering the best ranking quality (MRR = 0.627) while requiring less memory for smolagent multi-step reasoning. These top-performing models are also the fastest, taking approx. 21–22 seconds per query, indicating no latency penalty for improved accuracy. The BM25 baseline remains highly competitive (78% HR), outperforming most agent-based approaches and highlighting the continued strength of lexical matching in academic retrieval. Finally, RA-Bench results show higher performance than SemanticBench, suggesting that LLM-perturbed queries may be easier for multi-agent retrieval, though this requires further investigation.

Paper analysis visualization. In the Figure 3, we provide various output visualizations, including concept built graph (A), concept definition chart (B), interactive Q&A with precise information (C), markdown analysis output (D), and finally flow chart connecting the concepts of blocks (E). All of this analysis together provides the complete understanding of the paper.

Paper review analysis To evaluate our multi-agent review system, we conducted a study using the released ICLR 2024 reviews. We randomly selected 50 papers spanning diverse rating levels, and report the results in Figure 6. We observe that the code-oriented agent (qwen3-coder-30B) often struggles to sustain a coherent review workflow, whereas chat-style LLMs (e.g., gpt-oss) produce stronger and more consistent reviews. Overall, review quality improves with larger models, suggesting that capacity and instruction-following are particularly important for end-to-end reviewing.

4.3 Ablation Studies

We conduct comprehensive ablation studies to understand the contribution of different system com-

Table 2: Combined benchmark results for agent-based models and retrieval baselines. Best results are shown in **bold**. All the results are calculated using semantic benchmarks. Only the last (blue) is evaluated on 500 RABench queries, which shows syntetically written query is easier to retrieve compared to the random template following.

Model/Method	Type	Success	Hit Rate	MRR	R@1	R@5	R@10	R@20	R@50	Time (s)	Steps
Qwen3C-30B-Inst-Q3_K_M	Agent	100%	0.80	0.627	0.58	0.66	0.74	0.78	0.80	22.2	1.42
qwen3-coder:30b (Team, 2025)	Agent	100%	0.80	0.518	0.46	0.52	0.72	0.76	0.80	21.1	1.34
BM25 (Chen and Wiseman, 2023)	Baseline	100%	0.78	0.541	0.48	0.60	0.66	0.78	0.78	–	–
microcoder-deepseekr1-14.8	Agent	52%	0.73	0.453	0.38	0.46	0.65	0.69	0.73	107.4	4.15
deepseek-coder-v3:16b (Zhu et al., 2024a)	Agent	100%	0.66	0.396	0.32	0.46	0.52	0.60	0.66	47.9	1.54
qwen2.5-coder:3b (Hui et al., 2024)	Agent	94%	0.60	0.366	0.28	0.45	0.53	0.55	0.57	210.4	1.51
qwen2.5-coder:14b (Hui et al., 2024)	Agent	82%	0.56	0.461	0.41	0.51	0.51	0.56	0.56	73.4	1.05
Semantic (Wang et al., 2020)	Baseline	100%	0.54	0.279	0.22	0.32	0.38	0.52	0.54	–	–
Simple (bag-of-words)	Baseline	100%	0.54	0.279	0.22	0.32	0.38	0.52	0.54	–	–
qwen2.5-coder:7b (Hui et al., 2024)	Agent	100%	0.54	0.311	0.26	0.36	0.40	0.52	0.54	59.3	0.84
Qwen3C-30B-Inst-Q3_K_M	Agent	100%	0.42	0.348	0.32	0.38	0.38	0.40	0.42	22.7	1.40
deepseek-coder:33b (Zhu et al., 2024a)	Agent	100%	0.12	0.087	0.08	0.08	0.12	0.12	0.12	180.4	0.14
qwen3vl-4b-orlex	Agent	12%	0.08	0.080	0.08	0.08	0.08	0.08	0.08	37.9	0.14
granite-code:34b (Mishra et al., 2024)	Agent	100%	0.02	0.010	0.00	0.02	0.02	0.02	0.02	111.3	0.04
Hybrid (BM25+semantic)	Baseline	100%	0.02	0.001	0.00	0.00	0.00	0.00	0.02	–	–
qwen2.5-coder:1.5b (Hui et al., 2024)	Agent	100%	0.00	0.000	0.00	0.00	0.00	0.00	0.00	63.7	0.00
microcoder-oss-20b	Agent	54%	0.00	0.000	0.00	0.00	0.00	0.00	0.00	47.6	0.00
Qwen3-Coder-30B-A3B-Inst-Q3_K_M	Agent	100%	0.98	0.882	0.83	0.93	0.95	0.96	0.97	21.53	1.36

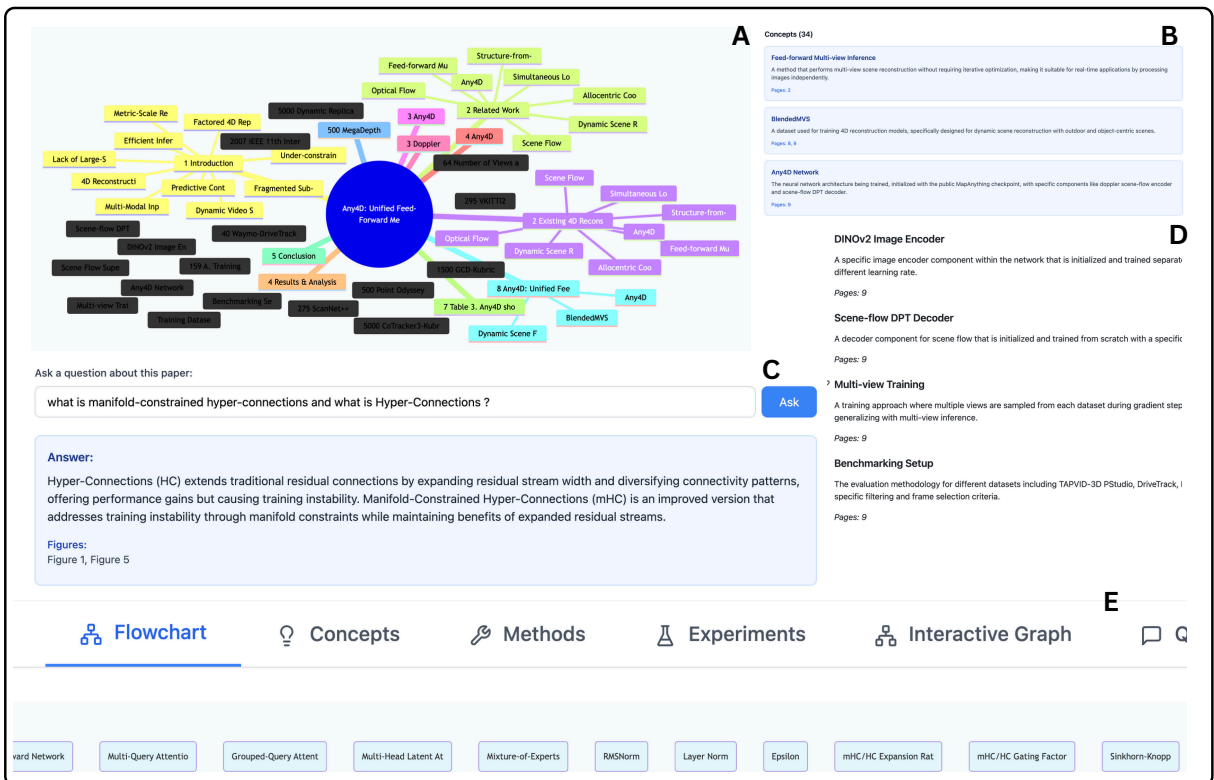


Figure 5: The main outputs of the analysis agent for a representative paper. (A) Interactive concept graph constructed from the paper, where nodes correspond to extracted concepts and edges denote semantic relationships. (B) Automatically generated concept explanations, each linked to the originating paper sections and pages. (C) Graph-aware question answering interface, providing answers grounded in extracted content along with supporting figures and references. (D) Structured Markdown exports summarizing all extracted concepts and methods for downstream use. (E) Flowchart view illustrating the high-level organization and relationships among concepts, methods, and experimental components of the paper.

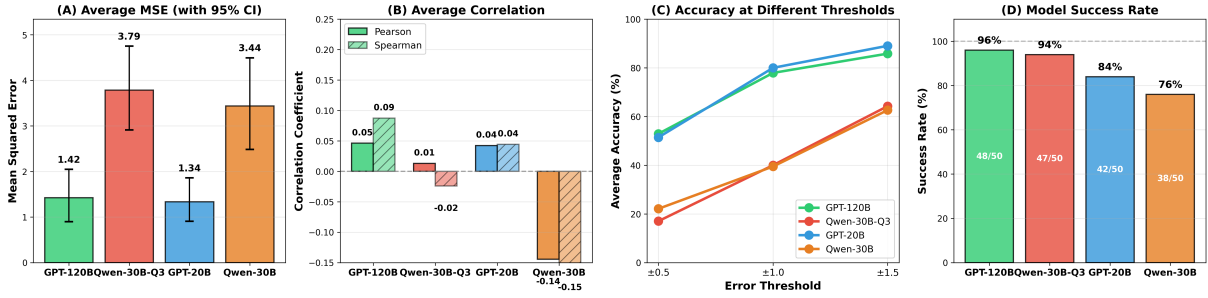


Figure 6: Paper review results analysis. This study was conducted on 50 randomly selected ICLR 2024 reviews.

ponents, including retrieval baselines, query configuration, and pipeline structures. **Full Query utilization**

To assess the full capability of our system, we conducted an extended evaluation using the qwen3-coder-30b model across 500 queries under various configurations. Results are presented in Table 3.

Table 3: Extended benchmark results for the Qooba agent (qwen3-coder-30b) across different configurations.

Configuration	Queries	Hit Rate	MRR	R@1	R@5	Time (s)
Default (Full Agent)	500	0.9818	0.8824	0.8381	0.9312	21.54
With Filters & Offline	50	0.9600	0.8485	0.7800	0.9000	22.76
Offline Only	50	0.9200	0.6476	0.5600	0.7400	41.45
No Mentions	50	0.6400	0.4316	0.3600	0.5200	38.35
Online/Offline Mix	50	0.6200	0.4595	0.4200	0.5000	38.50

Observations. The “With Filters & Offline” configuration performs better, suggesting that explicit context (conference/year filters) combined with local database access is highly effective. Notably, the “No Mentions” and “Online/Offline Mix” configurations show significant performance degradation (62–64% hit rate), indicating that specific paper references and structured retrieval chains are critical for accuracy. Overall, configurations exhibit similar latency, indicating stable MRR scaling of the multi-agent pipeline across query settings as well.

4.4 Retrieval Baseline Ablations

Table 4: Ablation study results comparing retrieval baselines and pipeline structures using qwen3-coder-30b. Full represents the full pipeline structure, minimal represents

Configuration	Baseline	Structure	Hit Rate	MRR	R@1	R@5	Time (s)
BM25 Full	bm25	full	0.9600	0.8629	0.8000	0.9200	33.75
BM25 Search Sort	bm25	search_sort	0.9600	0.8620	0.8000	0.9200	33.95
BM25 No Intent	bm25	no_intent	0.9600	0.8554	0.8000	0.9200	31.47
BM25 Search Analysis	bm25	search_analysis	0.9600	0.8437	0.7800	0.9200	32.81
BM25 Minimal	bm25	minimal	0.9600	0.8420	0.7800	0.9200	33.34
Hybrid Full	hybrid	full	0.9600	0.8620	0.8000	0.9200	31.65
BM25 + Reranker	bm25+reranker	full	0.9600	0.8692	0.8000	0.9400	935.07
Semantic Full	semantic	full	0.9400	0.7097	0.6200	0.8800	31.28

Retrieval Baseline Impact. BM25-based methods consistently outperform pure semantic retrieval. The semantic baseline shows a significant drop in R@1 (0.62) compared to BM25-based methods (0.80), suggesting that lexical matching remains crucial for precise paper retrieval. The hybrid approach performs on par with BM25, indicating that combining lexical and semantic signals does not provide additional benefits in this setting.

Reranking Trade-offs. The BM25 + Reranker configuration achieves the highest MRR (0.8692) and R@5 (0.9400), but at a substantial computational cost, approximately 28× slower than other methods. This presents a clear accuracy-efficiency trade-off that practitioners must consider based on their deployment requirements.

Pipeline Complexity. Reducing pipeline complexity (Minimal, Search Analysis configurations) leads to slight drops in MRR and R@1 while maintaining high overall hit rates (96%). Interestingly, removing intent analysis (“No Intent” configuration) results in a faster pipeline with competitive performance, suggesting that intent classification may be redundant for well-structured queries.

5 Conclusion

Paper Circle shows how multi-agent workflows can streamline research literature management. Its discovery pipeline unifies heterogeneous search sources and multi-criteria scoring into a reproducible tool, using a simple agent–tool interface with shared state, deterministic ranking, and synchronized multi-format outputs. Its analysis pipeline converts papers into structured knowledge graphs that enable graph-aware QA, coverage checks, and human-in-the-loop verification. Future work will focus on the optimization of the unification of the pipeline.

6 Limitations

Our review agent shows weak alignment with human judgments: across models, the correlation with human reviewer scores remains low ($r < 0.25$), and several metrics can even exhibit negative correlations, indicating that the system may rank papers in the opposite order of human preference. As a result, even the best-performing configurations do not reliably distinguish strong from weak submissions, and the system should not be used as a trusted mechanism for comparing or ranking papers. Based on our analysis, we found that this review process gets the benefit of a large model, so this problem can be overcome by large open/closed source models.

References

Narayan S Adhikari and Shradha Agarwal. 2024. A comparative study of pdf parsing tools across diverse document categories. *arXiv preprint arXiv:2410.09871*.

Prof. Chandrashekar, M. Akram, Mohin Khan, Piyush Kumar, and Pratap Mandal. 2025. *A survey on stock investment risk analysis using crewai multi-agent system*. *International Research Journal of Modernization in Engineering Technology and Science*.

Guoxin Chen, Zile Qiao, Wenqing Wang, Donglei Yu, Xuanzhong Chen, Hao Sun, Minpeng Liao, Kai Fan, Yong Jiang, Wayne Xin Zhao, and 1 others. 2025a. Mars: Optimizing dual-system deep research via multi-agent reinforcement learning. *arXiv preprint arXiv:2510.04935*.

Renqi Chen, Haoyang Su, SHIXIANG TANG, Zhenfei Yin, Qi Wu, Hui Li, Ye Sun, Wanli Ouyang, Philip Torr, and Nanqing Dong. 2025b. *Ai-driven automation can become the foundation of next-era science of science research*. *NIPS 2025*.

Xiaoyin Chen and Sam Wiseman. 2023. Bm25 query augmentation learned end-to-end. *arXiv preprint arXiv:2305.14087*.

KuanChao Chu, Yi-Pei Chen, and Hideki Nakayama. 2024. *Cohesive conversations: Enhancing authenticity in multi-agent simulated dialogues*. *COLM 2024*.

Mike D’Arcy, Tom Hope, Larry Birnbaum, and Doug Downey. 2024. Marg: Multi-agent review generation for scientific papers. *arXiv preprint arXiv:2401.04259*.

Mamata Das, PJA Alphonse, and 1 others. 2023. A comparative study on tf-idf feature weighting method and its analysis using unstructured dataset. *arXiv preprint arXiv:2308.04037*.

Anıl Dođru, Mehmet Onur Keskin, Catholijn M. Jonker, Tim Baarslag, and Reyhan Aydođan. 2024. *Negolog: An integrated python-based automated negotiation framework with enhanced assessment components*. *IJCAI 2024*.

Yao Fehlis, Charles Crain, Aidan Jensen, Michael Watson, James Juhasz, Paul Mandel, Betty Liu, Shawn Mahon, Daren Wilson, and Nick Lynch-Jonely. 2025. *Accelerating drug discovery through agentic ai: A multi-agent approach to laboratory automation in the dmta cycle*. *arXiv.org*.

Prof. Anjali Gajjar. 2025. *Evoresearch: A multi-agent ai framework for automated paper analysis*. *International Journal of Innovative Research in Advanced Engineering*.

Ruihao Gong, Shihao Bai, Siyu Wu, Yunqian Fan, Zaijun Wang, Xiuhong Li, Hailong Yang, and Xianglong Liu. 2025. Past-future scheduler for llm serving under sla guarantees. In *Proceedings of the 30th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*, pages 798–813.

Dong Han, Zhehong Ai, Pengxiang Cai, Shanya Lu, Jianpeng Chen, Zihao Ye, Shuzhou Sun, Ben Gao, Lingli Ge, Weida Wang, and 1 others. 2025a. Chembomas: Accelerated bo in chemistry with llm-enhanced multi-agent system. *arXiv preprint arXiv:2509.08736*.

Rujun Han, Yanfei Chen, Zoey CuiZhu, Lesly Miculicich, Guan Sun, Yuanjun Bi, Weiming Wen, Hui Wan, Chunfeng Wen, Solène Maître, and 1 others. 2025b. Deep researcher with test-time diffusion. *arXiv preprint arXiv:2507.16075*.

Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Kai Dang, and 1 others. 2024. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*.

Yiqiao Jin, Qinlin Zhao, Yiyang Wang, Hao Chen, Kaijie Zhu, Yijia Xiao, and Jindong Wang. 2024. Agentreview: Exploring peer review dynamics with llm agents. *arXiv preprint arXiv:2406.12708*.

Jiaxin Ju, YIZHEN ZHENG, Huan Yee Koh, Can Wang, and Shirui Pan. 2025. *Chemthinker: Thinking like a chemist with multi-agent llms for deep molecular insights*. *ICLR 2025*.

Mehmet Onur Keskin, Berk Buzcu, Berkecan Koçyiđit, Umut Çakan, Anıl Dođru, and Reyhan Aydođan. 2024. *Negotiator: A comprehensive framework for human-agent negotiation integrating preferences, interaction, and emotion*. *IJCAI 2024*.

Arpan Shaileshbhai Korat. 2025. *Synergistic minds: A collaborative multi-agent framework for integrated ai tool development using diverse large language models*. *World Journal of Advanced Research and Reviews*.

606	Shrinidhi Kumbhar, Venkatesh Mishra, Kevin Coutinho, Divij Handa, Ashif Iqbal, and Chitta Baral. 2025. Hypothesis generation for materials discovery and design using goal-driven and constraint-guided llm agents. <i>NAACL 2025</i> .	660
607		661
608		662
609		663
610		
611	Hao Duong Le, Xin Xia, and Chen Zhang. 2025. Multi-agent causal discovery using large language models. <i>ICLR 2025</i> .	664
612		665
613		666
614	Weizhen Li, Jianbo Lin, Zhuosong Jiang, Jingyi Cao, Xinpeng Liu, Jiayu Zhang, Zhenqiang Huang, Qianben Chen, Weichen Sun, Qiexiang Wang, and 1 others. 2025. Chain-of-agents: End-to-end agent foundation models via multi-agent distillation and agentic rl. <i>arXiv preprint arXiv:2508.13167</i> .	667
615		668
616		669
617		670
618		671
619		672
620	Mayank Mishra, Matt Stallone, Gaoyuan Zhang, Yikang Shen, Aditya Prasad, Adriana Meza Soria, Michele Merler, Parameswaran Selvam, Saptha Surendran, Shivdeep Singh, and 1 others. 2024. Granite code models: A family of open foundation models for code intelligence. <i>arXiv preprint arXiv:2405.04324</i> .	673
621		674
622		675
623		676
624		677
625		678
626	Vladimir Naumov, Diana Zagirova, Sha Lin, Yupeng Xie, Wenhao Gou, Anatoly Urban, Nina Tikhonova, Khadija M. Alawi, Mike Durymanov, and Fedor Galkin. 2025. Dora ai scientist: Multi-agent virtual research team for scientific exploration discovery and automated report generation. <i>bioRxiv</i> .	679
627		680
628		681
629		682
630		683
631		684
632	Seyednami Niyakan and Xiaoning Qian. 2025. Phenograph: A multi-agent framework for phenotype-driven discovery in spatial transcriptomics data augmented with knowledge graphs. <i>bioRxiv</i> .	685
633		686
634		687
635		688
636	Aniket Pramanick, Yufang Hou, Saif M. Mohammad, and Iryna Gurevych. 2023. A diachronic analysis of paradigm shifts in nlp research: When, how, and why? <i>EMNLP 2023</i> .	689
637		690
638		691
639		692
640	Yingming Pu, Tao Lin, and Hongyu Chen. 2025. Piflow: Principle-aware scientific discovery with multi-agent collaboration. <i>arXiv preprint arXiv:2505.15047</i> .	693
641		694
642		695
643		696
644	Renyi Qu, Ruixuan Tu, and Forrest Bao. 2025. Is semantic chunking worth the computational cost? In <i>Findings of the Association for Computational Linguistics: NAACL 2025</i> , pages 2155–2177.	697
645		698
646		699
647	Chandan K Reddy and Parshin Shojaee. 2025. Towards scientific discovery with generative ai: Progress, opportunities, and challenges. <i>AAAI 2025</i> .	700
648		701
649		702
650	Aymeric Roucher, Albert Villanova del Moral, Thomas Wolf, Leandro von Werra, and Erik Kaunismäki. 2025. ‘smolagents’: a smol library to build great agentic systems. https://github.com/huggingface/smolagents .	703
651		704
652		705
653		706
654		707
655	Alireza Salemi, Mihir Parmar, Palash Goyal, Yiwen Song, Jinsung Yoon, Hamed Zamani, Hamid Palangi, and Tomas Pfister. 2025. Llm-based multi-agent blackboard system for information discovery in data science. <i>arXiv preprint arXiv:2510.01285</i> .	708
656		709
657		710
658		711
659		712
		713
	Er. Jagpreet Singh and Prasant Kumar. 2025. <i>Astrafin: ai financial agent</i> . <i>INTERNATIONAL JOURNAL OF SCIENTIFIC RESEARCH IN ENGINEERING AND MANAGEMENT</i> .	
	Jackson Spieser, Ali Balapour, Jarek Meller, Krushna Patra, and Behrouz Shamsaei. 2025. Multi-agent ai systems for biological and clinical data analysis. <i>Preprints.org</i> .	
	Qwen Team. 2025. <i>Qwen3 technical report</i> . <i>Preprint</i> , arXiv:2505.09388.	
	Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. <i>Advances in neural information processing systems</i> , 33:5776–5788.	
	Xinyue Wang, Kun Zhou, Wenyi Wu, Har Simrat Singh, Fang Nan, Songyao Jin, Aryan Philip, Saloni Patnaik, Hou Zhu, Shivam Singh, and 1 others. 2025. Causal-copilot: An autonomous causal analysis agent. <i>arXiv preprint arXiv:2504.13263</i> .	
	Michał Wawer and Jarosław A. Chudziak. 2025. Integrating traditional technical analysis with ai: A multi-agent llm-based approach to stock market forecasting. <i>International Conference on Agents and Artificial Intelligence</i> .	
	Michael Wooldridge. 2002. <i>An introduction to multi-agent systems</i> . John Wiley & Sons.	
	Mengxi Xiao, Ben Liu, He Li, Jimin Huang, Qianqian Xie, Xiaofen Zong, Mang Ye, and Min Peng. 2025. Moodangels: A retrieval-augmented multi-agent framework for psychiatry diagnosis. <i>NIPS 2025</i> .	
	Yihang Xiao, Jinyi Liu, Yan Zheng, Xiaohan Xie, Jianye Hao, Mingzhi Li, Ruitao Wang, Fei Ni, Yuxiao Li, Jintian Luo, and 1 others. 2024. Cellagent: An llm-driven multi-agent framework for automated single-cell data analysis. <i>arXiv preprint arXiv:2407.09811</i> .	
	Zonglin Yang, Wanhao Liu, Ben Gao, Tong Xie, Yuqiang Li, Wanli Ouyang, Soujanya Poria, Erik Cambria, and Dongzhan Zhou. 2025. <i>Moose-chem: Large language models for rediscovering unseen chemistry scientific hypotheses</i> . <i>ICLR 2025</i> .	
	Haofei Yu, Zirui Cheng, Zhaochen Hong, Kunlun Zhu, Jinwei Yao, Tao Feng, and Jiakuan You. 2025a. <i>Research town: Simulator of research community</i> . <i>ICLR 2025</i> .	
	Haofei Yu, Zhaochen Hong, Zirui Cheng, Kunlun Zhu, Keyang Xuan, Jinwei Yao, Tao Feng, and Jiakuan You. 2025b. <i>Researchtown: Simulator of human research community</i> . <i>ICML 2025</i> .	
	Zhaojian Yu, Kaiyue Feng, Yilun Zhao, Shilin He, Xiaoping Zhang, and Arman Cohan. 2025c. <i>Alpharesearch: Accelerating new algorithm discovery with language models</i> . <i>arXiv preprint arXiv:2511.08522</i> .	

714 Yurun Yuan and Tengyang Xie. 2025. Reinforce llm rea- 765
715 soning through multi-agent reflection. *arXiv preprint* 766
716 *arXiv:2506.08379*. 767

717 Bohui Zhang, Yuan He, Lydia Pintscher, Albert Meroño 768
718 Peñuela, and Elena Simperl. 2025a. Schema genera- 769
719 tion for large knowledge graphs using large language 770
720 models. *arXiv preprint arXiv:2506.04512*. 771

721 Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, 772
722 Huan Lin, Baosong Yang, Pengjun Xie, An Yang, 773
723 Dayiheng Liu, Junyang Lin, Fei Huang, and Jingren 774
724 Zhou. 2025b. Qwen3 embedding: Advancing text 775
725 embedding and reranking through foundation models. 776
726 *arXiv preprint arXiv:2506.05176*. 777

727 Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, 778
728 Peiyi Wang, Runxin Xu, Y Wu, Yukun Li, Huazuo 779
729 Gao, Shirong Ma, and 1 others. 2024a. Deepseek- 780
730 coder-v2: Breaking the barrier of closed-source 781
731 models in code intelligence. *arXiv preprint* 782
732 *arXiv:2406.11931*. 783

733 Yuqi Zhu, Xiaohan Wang, Jing Chen, Shuofei Qiao, 784
734 Yixin Ou, Yunzhi Yao, Shumin Deng, Huajun Chen, 785
735 and Ningyu Zhang. 2024b. LLMs for knowledge 786
736 graph construction and reasoning: Recent capabilities 787
737 and future opportunities. *World Wide Web*, 27(5):58. 788

738 A Paper Review Results

739 We evaluate how well large language models 790
740 can predict human paper-review scores on ICLR 791
741 submissions. From the ICLR 2024 dataset, 792
742 we randomly sampled 50 papers to cover a 793
743 broad range of human-assigned ratings and eval- 794
744 uated four tool-enabled LLMs: gpt-oss:120b, 795
745 gpt-oss:20b, qwen3-coder-30b, and a quantized 796
746 qwen3-coder-30b variant. For each paper, the 797
747 model produces numerical scores for standard re- 798
748 view dimensions (overall rating, soundness, pre- 799
749 sentation, and contribution), which we compare 800
750 against the corresponding human scores. 801

751 **Metrics.** We report regression error (MSE, MAE, 802
752 RMSE), rank/linear association (Pearson, Spear- 803
753 man), and thresholded accuracy (percentage of pre- 804
754 dictions within ± 0.5 , ± 1.0 , and ± 1.5 of the human 805
755 score). We also report the mean and standard devia- 806
756 tion of signed errors to characterize systematic bias. 807
757 Due to occasional missing fields or filtering during 808
758 preprocessing, the number of evaluated papers N 809
759 can differ slightly across models. 810

760 **Key findings.** Across categories, gpt-oss:120b 811
761 achieves the best overall accuracy on *rating* and 812
762 *contribution* (e.g., rating MAE = 1.68; contribu- 813
763 tion MAE = 0.62), while gpt-oss:20b is com-
764 petitive and often stronger on more technical sub-

scores such as *soundness* and *presentation*. De-
spite moderate absolute errors on several dimen-
sions, correlations with human scores remain weak
across models (generally $|r| < 0.25$), suggesting
that models struggle to preserve the relative rank-
ing of papers even when their average deviation is
limited. Code-specialized models (Qwen3-Coder)
remain viable baselines, but show larger errors on
overall rating and contribution in this setting.

774 B System Overview

775 Paper Circle is a full-stack platform with a web fron- 776
777 tend and a Python backend as shown in the Figure 7. 778
779 The frontend (React, TypeScript, Vite, TailwindCSS) 780
781 provides discovery, reading circles, and discussion 782
782 features. The backend exposes discovery APIs via 783
783 FastAPI and implements the multi-agent pipelines 784
784 used by the system. Supabase (PostgreSQL + Auth) 785
785 provides storage for users, communities, papers, and 786
786 sessions. 787

788 The discovery backend includes two major 789
789 pipelines: (i) a refactored research discovery 790
790 pipeline focused on deterministic retrieval, scor- 791
791 ing, and diversity, and (ii) a multi-agent research 792
792 pipeline that produces structured step-by-step out- 793
793 puts with offline search support. Both pipelines 794
794 are accessible through API endpoints and are in- 795
795 tegrated into the Paper Circle user interface for 796
796 interactive discovery workflows. 797

793 B.1 System Architecture

794 Figure 8 illustrates the overall architecture of Paper 795
795 Circle. The system consists of two complementary 796
796 multi-agent pipelines: the *Discovery Pipeline* for 797
797 finding relevant papers, and the *Analysis Pipeline* 798
798 for deep understanding of individual papers. 799

799 The discovery pipeline, as shown in the Fig- 800
800 ure 8 is composed of six agents: intent classifica- 801
801 tion, paper search, sorting, analysis, export, and 802
802 web search. The intent classifier parses natural- 803
803 language queries into structured constraints (search 804
804 mode, conferences, year range, max results, and 805
805 ranking preferences). The paper search agent is the 806
806 primary retrieval worker; it updates the global state 807
807 and writes outputs after every search step. The 808
808 sorting and analysis agents operate on the shared 809
809 paper list to refine ranking and derive insights. The 810
810 export agent centralizes output access for down- 811
811 stream workflows, while the web search agent sup- 812
812 plements the pipeline with external lookup tools 813
813 when required. All agents are coordinated by the

Model N	Category	MSE	MAE	RMSE	Pearson	Spearman	Acc. ± 0.5	Acc. ± 1.0	Acc. ± 1.5	Mean Err.	Std Err.
oss-120B 48	RATING	4.6934	1.6844	2.1664	-0.0407	0.0571	25.00%	43.75%	58.33%	0.2177	2.1555
oss-120B 48	SOUNDNESS	0.7316	0.6351	0.8554	-0.0054	0.0474	58.33%	85.42%	87.50%	-0.0816	0.8515
oss-120B 48	PRESENTATION	0.6564	0.6038	0.8102	0.0701	0.1259	60.42%	83.33%	91.67%	-0.0920	0.8049
oss-120B 48	CONTRIBUTION	0.6349	0.6240	0.7968	0.0717	0.0734	56.25%	85.42%	91.67%	0.0087	0.7967
oss-20 42	RATING	4.7607	1.7647	2.1819	0.0989	0.1869	21.43%	40.48%	52.38%	1.5980	1.4856
oss-20 42	SOUNDNESS	0.4241	0.5190	0.6512	-0.0106	-0.0226	59.52%	92.86%	97.62%	0.3294	0.5618
oss-20 42	PRESENTATION	0.4271	0.5171	0.6535	-0.1270	-0.1299	64.29%	90.48%	97.62%	0.3512	0.5511
oss-20 42	CONTRIBUTION	0.6482	0.6702	0.8051	0.2221	0.1757	50.00%	83.33%	97.62%	0.6250	0.5075
qwen30B-code_qk_3 47	RATING	11.8533	2.9879	3.4429	-0.2233	-0.2837	8.51%	17.02%	29.79%	2.9085	1.8422
qwen30B-code_qk_3 47	SOUNDNESS	1.6941	1.1730	1.3016	0.0113	-0.0096	17.02%	46.81%	72.34%	1.1454	0.6182
qwen30B-code_qk_3 47	PRESENTATION	1.4257	1.0191	1.1940	0.0378	0.0271	27.66%	59.57%	78.72%	0.9787	0.6840
qwen30B-code_qk_3 47	CONTRIBUTION	2.2921	1.3865	1.5140	0.0196	0.0224	12.77%	34.04%	65.96%	1.3865	0.6080
Qwen 30B 38	RATING	10.2331	2.7930	3.1989	-0.1820	-0.2216	7.89%	13.16%	26.32%	2.6930	1.7266
Qwen 30B 38	SOUNDNESS	1.7172	1.2096	1.3104	-0.1157	-0.1057	13.16%	39.47%	73.68%	1.1491	0.6298
Qwen 30B 38	PRESENTATION	0.9526	0.7180	0.9760	-0.1319	-0.1495	55.26%	73.68%	81.58%	0.6522	0.7261
Qwen 30B 38	CONTRIBUTION	2.5212	1.4746	1.5878	-0.2119	-0.2160	13.16%	26.32%	55.26%	1.4640	0.6146

Table 5: **Paper review score prediction on ICLR 2024.** We compare four LLMs on predicting human review scores across rating, soundness, presentation, and contribution. We report error metrics (MSE/MAE/RMSE), correlation (Pearson/Spearman), and thresholded accuracy (within ± 0.5 , ± 1.0 , ± 1.5 of the human score). N denotes the number of papers evaluated for each model after preprocessing.

CodeAgent, which enforces a minimal-step policy for efficiency and uses the intent classifier to decide offline versus online search.

The analysis pipeline operates on individual papers, transforming PDF documents into structured knowledge graphs. It employs four specialized extraction agents (concept, method, experiment, and linkage) that process paper content in phases, building a typed graph with full traceability to source locations. The resulting graph supports question answering, coverage verification, and multi-format export.

B.2 State Management and Outputs

State is maintained in PipelineState. Each step increments a counter, logs action metadata, and regenerates synchronized artifacts. The outputs include: (i) papers.json with full paper metadata and computed scores, (ii) links.json with structured links and PDFs/DOIs, (iii) stats.json with aggregate statistics and a leaderboard, (iv) summary.json with generated insights and key findings, (v) retrieval_metrics.json when evaluation is enabled, and (vi) human-readable exports (CSV, BibTeX, Markdown) plus a live HTML

dashboard. This approach ensures that each agent step is reproducible and auditable.

B.3 Retrieval

The pipeline supports both offline and online retrieval. Offline search loads papers from a local JSON corpus and optionally filters by conference and year. It ranks results using BM25 by default, with optional semantic similarity (sentence transformers) or hybrid scoring when available. An optional cross-encoder reranker can refine the top results; when enabled, it reranks a first-stage candidate set. Online search aggregates results from arXiv, Semantic Scholar, OpenAlex, and DBLP via their public APIs. A query intent classifier detects search mode, conference constraints, year ranges, and ranking preferences, and routes the query to the appropriate retrieval pathway. Deduplication is applied across sources by normalizing titles.

B.4 Ranking and Scoring

After retrieval, papers are scored along multiple axes: recency, similarity to the query (TF-IDF (Das et al., 2023) when available), novelty based on title token frequency, and normalized BM25 scores (Chen and Wiseman, 2023). The

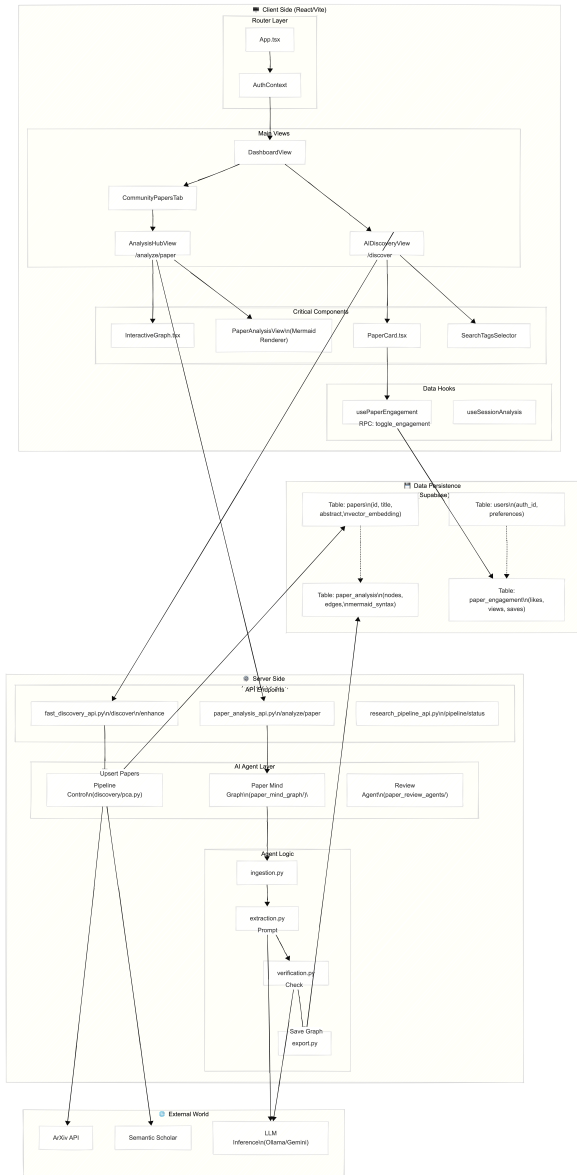


Figure 7: Frontend of the website from client side.

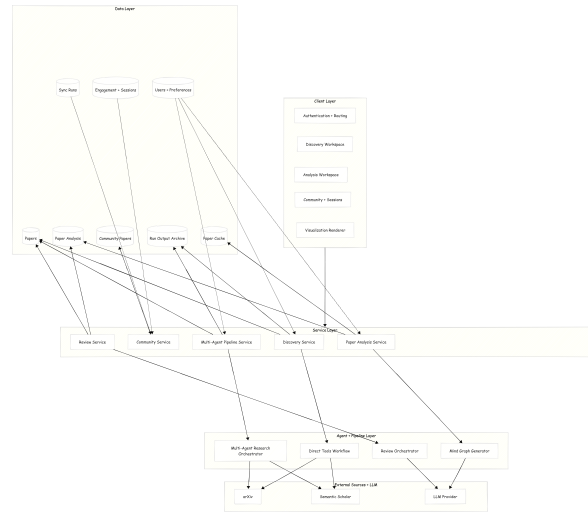


Figure 8: Discovery pipeline front-end view of the website, which takes an input query and iteratively refines the state and updates the database.

system supports sorting by any single criterion or by a weighted combined score. Relevance scores are computed as a weighted mixture of similarity, recency, citation count, and BM25. Final ranks are assigned after sorting, and the updated ordering is reflected in all exported artifacts. When reranker-based sorting is requested, a cross-encoder replaces the default scoring with direct relevance scores.

B.5 Analysis and Monitoring

The pipeline computes aggregate statistics such as source distribution, year distribution, top authors and venues, keyword frequency, and citation summaries. These analytics populate structured summaries and are visualized in an auto-refreshing HTML dashboard. Each agent action is logged with timestamps and paper counts, enabling reproducibility and step-level auditing of the pipeline. The pipeline also maintains a step log that captures the agent name, action, results preview, and parameters used.

C Retrieval Pipeline

Paper Circle supports both offline and online retrieval to balance coverage, speed, and reproducibility. The choice between retrieval modes is controlled by the intent classification agent, which parses user queries to determine the optimal search strategy.

C.1 Offline Retrieval

The OfflinePaperSearchEngine enables fast (See the Figure 9, reproducible search over a local

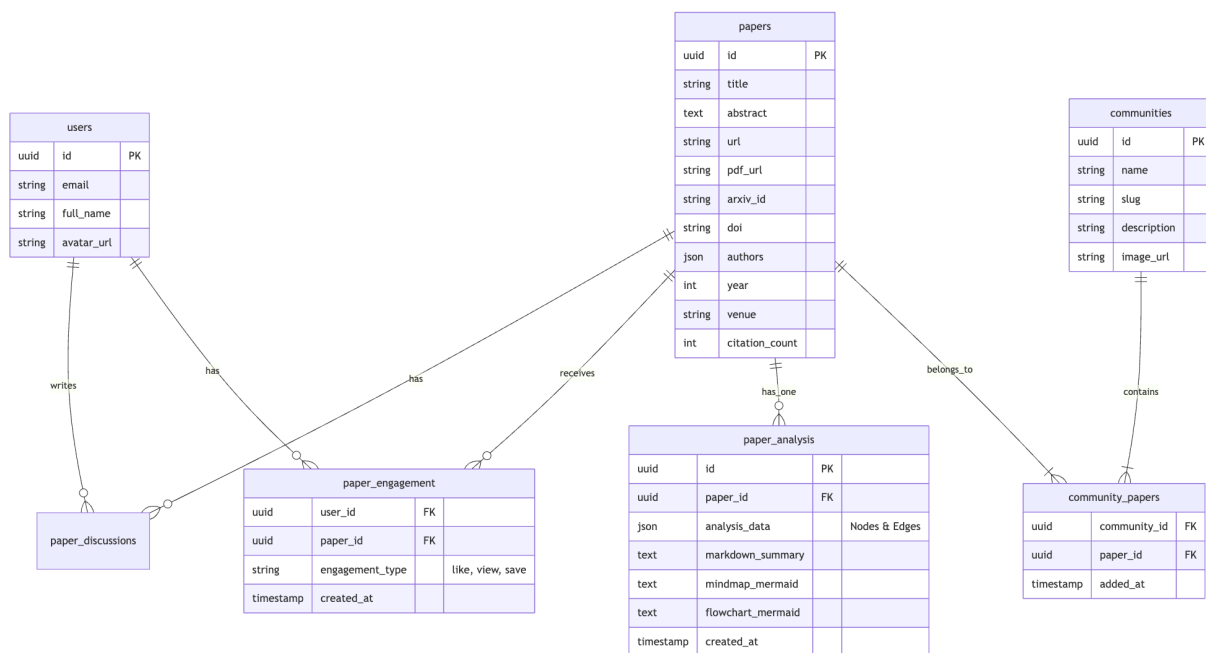


Figure 9: Paper analysis and database management for fast inference.

892 database of academic papers stored as JSON files.
 893 Each database file contains structured paper meta-
 894 data including title, authors, abstract, venue, year,
 895 track, keywords, and DOI.

896 The offline search process:

- 897 1. **Database Loading:** Papers are loaded from
 898 the specified database path with optional fil-
 899 tering by conference (e.g., ICLR, NeurIPS,
 900 ACL) and year range.
- 901 2. **Text Preparation:** For each paper, searchable
 902 text is constructed by concatenating the title,
 903 abstract, and keywords.
- 904 3. **BM25 Indexing:** When available, papers are
 905 indexed using the Okapi BM25 algorithm via
 906 the rank_bm25 library. The index uses tok-
 907 enized documents for sparse retrieval.
- 908 4. **Query Execution:** User queries are tokenized
 909 and scored against the BM25 index, returning
 910 a ranked list of candidates.

911 An optional cross-encoder reranker can refine
 912 the top-*k* results from the first-stage retrieval.
 913 When enabled via the AdvancedReranker mod-
 914 ule, the system uses a transformer-based reranker
 915 (e.g., Qwen3-Reranker) to compute more precise
 916 relevance scores between the query and candidate
 917 documents.

918 C.2 Online Retrieval

919 For broader or more current searches, Paper Circle
 920 aggregates results from multiple academic APIs:

- 921 • **arXiv:** Queries the arXiv API for preprints,
 922 extracting title, authors, abstract, categories,
 923 and PDF links.
- 924 • **Semantic Scholar:** Retrieves papers with cita-
 925 tion counts, abstracts, and venue information
 926 via the Semantic Scholar Academic Graph
 927 API.
- 928 • **OpenAlex:** Accesses the OpenAlex catalog
 929 for open-access metadata and citation net-
 930 works.
- 931 • **DBLP:** Searches the DBLP computer science
 932 bibliography for venue-specific results.

933 Each source is queried in parallel using a thread
 934 pool executor for efficiency. Results are normal-
 935 ized into the common Paper data structure before
 936 merging.

937 C.3 Deduplication

938 After retrieval, the pipeline performs two-stage
 939 deduplication to eliminate redundant entries:

- 940 1. **DOI-based deduplication:** Papers with
 941 matching DOIs are deduplicated, preferring
 942 entries with richer metadata (e.g., abstracts,
 943 PDF URLs).

- 944 2. **Title-based deduplication:** Titles are normal- 985
 945 ized by removing punctuation and convert- 986
 946 ing to lowercase. Duplicate titles are merged,
 947 again preferring metadata-complete entries.

948 The deduplication step is critical when aggregating 988
 949 results from multiple sources, as the same 989
 950 paper often appears in arXiv, Semantic Scholar,
 951 and OpenAlex with varying metadata quality.

952 C.4 Query Expansion

953 The query generation agent converts natural-
 954 language user input into a structured search specifi-
 955 cation containing:

- 956 • **Core keywords:** Primary search terms extracted 995
 957 from the query. 996
- 958 • **Required constraints:** Mandatory terms that 997
 959 must appear in results. 998
- 960 • **Related terms:** Synonyms or related concepts 999
 961 to expand recall. 1000
- 962 • **Negative keywords:** Terms to exclude from 1001
 963 results. 1002
- 964 • **Plausible paper titles:** Hypothesized titles 1003
 965 for targeted retrieval. 1004

966 This structured specification enables consistent 1005
 967 query construction across heterogeneous data 1006
 968 sources while capturing user intent more precisely 1007
 969 than raw keyword matching. 1008

970 D Scoring and Ranking

971 Paper Circle employs a multi-criteria scoring frame- 1009
 972 work designed for research discovery rather than 1010
 973 general information retrieval. Each paper receives 1011
 974 scores along multiple dimensions, which are com- 1012
 975 bined using mode-specific weights to produce a 1013
 976 final ranking. 1014

977 D.1 Scoring Dimensions

978 The system computes the following scores for each 1015
 979 retrieved paper: 1016

980 **Similarity Score** Relevance to the user query is 1017
 981 computed using TF-IDF (Das et al., 2023) vector- 1018
 982 ization and cosine similarity. The query and pa- 1019
 983 per text (concatenated title and abstract) are trans- 1020
 984 formed into TF-IDF vectors using scikit-learn’s

TfidfVectorizer. The similarity score is the co- 985
 sine of the angle between these vectors: 986

$$987 \text{similarity}(q, p) = \frac{\vec{v}_q \cdot \vec{v}_p}{\|\vec{v}_q\| \cdot \|\vec{v}_p\|} \quad (1)$$

where \vec{v}_q and \vec{v}_p are the TF-IDF vectors for the 988
 query and paper, respectively. 989

Recency Score Papers are scored by publication 990
 year, with more recent papers receiving higher 991
 scores. The recency score is normalized relative to 992
 the current year: 993

$$994 \text{recency}(p) = \frac{\text{year}(p) - \text{year}_{\min}}{\text{year}_{\max} - \text{year}_{\min}} \quad (2)$$

where year_{\min} and year_{\max} are the minimum and 995
 maximum years in the corpus. 996

Novelty Score Novelty measures how different 997
 a paper is from the corpus centroid, computed as 998
 the TF-IDF distance from the average document 999
 vector. Papers with unusual terminology or unique 1000
 topic combinations receive higher novelty scores, 1001
 surfacing potentially overlooked works. 1002

BM25 Score When the rank_bm25 library is 1003
 available, the Okapi BM25 algorithm provides an 1004
 alternative relevance measure that accounts for 1005
 term frequency saturation and document length 1006
 normalization. BM25 scores are normalized to 1007
 the $[0, 1]$ range for comparability with other dimen- 1008
 sions. 1009

Citation Count When available from the source 1010
 API (primarily Semantic Scholar and OpenAlex), 1011
 citation counts provide a proxy for impact. Citation- 1012
 based ranking is optional and disabled by default 1013
 to avoid recency bias against new papers. 1014

1015 D.2 Combined Score Computation

The final combined score is a weighted sum of 1016
 individual dimensions: 1017

$$1018 \text{combined}(p) = w_s \cdot \text{similarity} + w_r \cdot \text{recency} + w_n \cdot \text{novelty} + w_b \cdot \text{bm25} \quad (3)$$

The weights (w_s, w_r, w_n, w_b) are determined by 1019
 the search mode: 1020

- 1021 • **Stable mode:** Prioritizes relevance and author- 1021
 1022 ity. Weights: $w_s = 0.5, w_r = 0.2, w_n = 0.1,$
 1023 $w_b = 0.2.$
- 1024 • **Discovery mode:** Prioritizes novelty to sur- 1024
 1025 face non-obvious results. Weights: $w_s = 0.3,$
 1026 $w_r = 0.1, w_n = 0.4, w_b = 0.2.$

- **Balanced mode:** Equal emphasis across dimensions. Weights: $w_s = 0.3$, $w_r = 0.2$, $w_n = 0.2$, $w_b = 0.3$.

Users can override these weights at query time via API parameters, enabling custom relevance trade-offs for specific research contexts.

D.3 Sorting Stage

After scoring, the sorting agent reorders papers according to user preferences. Supported sort criteria include:

- recency: Most recent papers first.
- citations: Highest-cited papers first.
- similarity: Most relevant papers first.
- novelty: Most unusual papers first.
- bm25: Best BM25 matches first.
- combined: Weighted combined score (default).

D.4 Cross-Encoder Reranking

For high-precision use cases, the pipeline supports optional cross-encoder reranking. When enabled, a transformer-based reranker (configured via `RerankerConfig`) processes query-document pairs through a cross-attention model to compute more accurate relevance scores than first-stage retrieval alone. The `MultiStageRetriever` first retrieves a larger candidate set (e.g., top-200) using BM25, then reranks to produce the final top- k results. This two-stage approach balances efficiency with ranking quality.

E Diversity and Postprocessing

Relevance-based ranking alone can produce homogeneous results, with multiple papers covering similar topics or methods. Paper Circle addresses this through diversity-aware postprocessing that ensures the top results span a broader range of perspectives.

E.1 Maximal Marginal Relevance

To improve topical coverage, Paper Circle applies Maximal Marginal Relevance (MMR) to the candidate list after initial scoring. MMR iteratively selects papers that maximize a combination of relevance to the query and dissimilarity to already-selected papers:

$$\text{MMR} = \arg \max_{p \in R \setminus S} \left[\lambda \cdot \text{sim}(p, q) - (1 - \lambda) \cdot \max_{s \in S} \text{sim}(p, s) \right] \quad (4)$$

where R is the candidate set, S is the set of already-selected papers, q is the query, and λ controls the relevance–diversity trade-off.

The diversity parameter λ is mode-dependent:

- **Stable mode:** $\lambda = 0.8$ (relevance-focused).
- **Discovery mode:** $\lambda = 0.5$ (diversity-focused).
- **Balanced mode:** $\lambda = 0.65$.

Similarity between papers is computed using TF–IDF cosine similarity over concatenated title and abstract text. This ensures that top results cover distinct subtopics rather than repeating variations of the same idea.

E.2 Secondary Views

The pipeline constructs specialized views over the ranked list to serve different discovery goals:

Hidden Gems Papers with high novelty scores but moderate relevance scores are surfaced as “hidden gems.” These are papers that may not rank highly on traditional relevance metrics but offer unique perspectives or cover underexplored topics. The hidden gems view is computed by sorting papers by novelty score and filtering for those below rank 20 in the combined ranking.

Canonical Papers Papers with high citation counts or appearing in top-tier venues are flagged as “canonical” works. This view helps users identify foundational papers in a research area, complementing the recency-focused main ranking.

Source Distribution The postprocessing stage also reports the distribution of papers across sources (arXiv, Semantic Scholar, etc.), enabling users to assess coverage and identify potential gaps in the retrieval.

E.3 Statistics and Analytics

After ranking, the analysis agent computes aggregate statistics stored in `stats.json`:

- **Year distribution:** Paper counts by publication year.
- **Source distribution:** Paper counts by retrieval source.

1112	• Top authors: Authors appearing most frequently in results.	• <code>stats.json</code> : Aggregate statistics and leaderboards.	1151
1113			1152
1114	• Top venues: Conferences and journals with highest representation.	• <code>summary.json</code> : Insights and key findings.	1153
1115		• <code>retrieval_metrics.json</code> : Step-level evaluation metrics.	1154
1116	• Keyword frequency: Most common terms in paper titles.		1155
1117		Additional exports include CSV, BibTeX, Markdown, and an auto-refreshing HTML dashboard.	1156
1118	• Citation statistics: Total, average, median, min, and max citation counts.	These outputs allow the same discovery session to be used for curation, citation management, and reporting.	1157
1119		The system exposes REST APIs via FastAPI. The discovery endpoint accepts a query and mode, returns structured search specifications, and provides the full ranked list with scores. Mode weights can be queried or overridden at runtime, enabling customized relevance/authority/novelty trade-offs.	1158
1120	• Score statistics: Average similarity, novelty, recency, and BM25 scores.		1159
1121			1160
1122	These analytics are visualized in an auto-refreshing HTML dashboard that updates every 10 seconds during pipeline execution, providing real-time visibility into the discovery process.		1161
1123			1162
1124			1163
1125			1164
1126	E.4 Insight Generation	G Evaluation	1167
1127	The pipeline automatically generates human-readable insights from the collected data:	We evaluate Paper Circle along three axes: (i) retrieval effectiveness under different configurations, (ii) stability and reproducibility of rankings across steps, and (iii) the utility of diversity-aware postprocessing for surfacing non-redundant results. Paper Circle provides built-in evaluation metrics but does not enforce a fixed benchmark dataset. When a ground-truth paper title or identifier is provided, the system computes Mean Reciprocal Rank (MRR), Recall@K, Precision@K, and hit rates. These metrics are computed per step and stored in JSON file for longitudinal tracking.	1168
1128		As a minimal illustrative scenario, consider a known target paper in the local corpus: the pipeline is run once using offline retrieval and once using online sources. The resulting MRR and Recall@K values allow direct comparison of configuration impact, while repeated runs confirm stable rankings when deterministic scoring is enabled. Although lightweight, this framing aligns evaluation with discovery goals rather than task-specific QA benchmarks.	1169
1129	• Publication trends: Identifies the year with the most publications.		1170
1130			1171
1131	• Primary source: Reports which API contributed the most results.		1172
1132			1173
1133	• Prolific authors: Highlights researchers with multiple papers in the collection.		1174
1134			1175
1135	• Citation leaders: Identifies the most-cited paper.		1176
1136			1177
1137	• Hot topics: Lists the most frequent keywords.		1178
1138	• Open access availability: Reports the percentage of papers with direct PDF links.		1179
1139			1180
1140	These insights are stored in <code>summary.json</code> and displayed on the dashboard, helping users quickly understand the landscape of retrieved literature.		1181
1141			1182
1142			1183
1143	F Outputs and Interfaces	For batch evaluation, a parallel benchmarking utility executes multiple queries concurrently and aggregates mean metrics and timing statistics. This supports lightweight comparisons between search configurations (offline vs. online, BM25 vs. semantic, with or without reranking) without requiring external tooling.	1184
1144	The pipeline maintains synchronized structured outputs after every agent step. The primary artifacts include:		1185
1145			1186
1146			1187
1147	• <code>papers.json</code> : Full paper metadata and scores.		1188
1148			1189
1149	• <code>links.json</code> : Structured links and PDF/DOI entries.	Knowledge Graph Schema. The mind graph follows a typed schema with nodes for papers, sections, concepts, methods, experiments, datasets,	1190
1150			1191
			1192
			1193
			1194
			1195
			1196
			1197
			1198
			1199

1200	and visual elements (figures, tables, equations), and	1. Retrieves the top- k most similar chunks and	1247
1201	edges encoding structural and semantic relations	nodes.	1248
1202	such as hierarchy, definition, proposal, usage, eval-	2. Expands context by including 1-hop graph	1249
1203	uation, illustration, and dependency. Each node	neighbors.	1250
1204	and edge is annotated with provenance metadata,	3. Returns chunks, nodes, and connecting edges.	1251
1205	including source chunk IDs, page numbers, verifi-		
1206	cation status, confidence scores, and timestamps,	The PaperQA agent constructs a prompt with	1252
1207	providing full traceability from any graph element	the retrieved context, including text chunks with	1253
1208	back to the original PDF.	their section sources, relevant concept descriptions,	1254
		and graph relationships. The response includes the	1255
1209	G.1 Multi-Agent Extraction	answer, supporting sections, relevant figures and	1256
1210	The GraphBuilder orchestrates four special-	tables, and a confidence estimate.	1257
1211	ized extraction agents, each implemented as a	A locate function allows users to find where	1258
1212	CodeAgent with domain-specific instructions:	specific items are discussed in the paper by search-	1259
1213	Concept Extractor Identifies key concepts from	ing across nodes, figures, tables, and text chunks,	1260
1214	text chunks, classifying each by type (definition,	returning page numbers and context snippets.	1261
1215	technique, theory, phenomenon) and importance		
1216	(core, supporting, background). The agent outputs	G.3 Coverage Verification	1262
1217	structured JSON with concept names, descriptions,	To ensure nothing is silently dropped during ex-	1263
1218	and classifications.	traction, the CoverageChecker produces a detailed	1264
		coverage report:	1265
1219	Method Extractor Focuses on sections con-	• Figure coverage: How many figures are	1266
1220	taining method-related keywords (“method”, “ap-	linked to concepts or methods.	1267
1221	proach”, “architecture”, “algorithm”). For each	• Table coverage: How many tables are linked	1268
1222	method, it extracts the name, description, category	to results or experiments.	1269
1223	(proposed, baseline, component), and key steps.	• Section coverage: How many sections have	1270
		extracted concepts.	1271
1224	Experiment Extractor Processes experiment	• Equation coverage: How many equations are	1272
1225	sections to extract experimental setups, datasets	linked to concepts they define.	1273
1226	used, evaluation metrics, and key results. It also	The report includes an overall coverage score	1274
1227	identifies dataset nodes for cross-referencing.	(0–100%), lists of unlinked items with suggestions,	1275
		and critical issues (e.g., “No figures are linked to	1276
1228	Linkage Agent Connects figures and tables to	concepts/methods”). This enables quality assur-	1277
1229	the concepts and methods they illustrate. Given	ance before downstream use.	1278
1230	a figure caption, nearby text, and a list of exist-	G.4 Human Verification Workflow	1279
1231	ing concepts, the agent determines which concepts	The VerificationManager supports human-in-	1280
1232	the figure relates to and the type of relationship	the-loop review:	1281
1233	(illustrates, summarizes, compares, demonstrates).	• <code>verify_node</code> : Mark a node as human-	1282
1234	The extraction proceeds in five phases: (1) con-	verified.	1283
1235	cept extraction from body chunks, (2) method ex-	• <code>edit_node</code> : Modify node title or description.	1284
1236	traction from method sections, (3) experiment and	• <code>add_edge</code> : Create new relationships.	1285
1237	dataset extraction, (4) figure and table linkage,	• <code>remove_edge</code> : Delete incorrect relationships.	1286
1238	and (5) inter-concept relationship discovery. Each	• <code>flag_for_review</code> : Flag nodes for review	1287
1239	phase updates the shared MindGraph data structure.	with a reason.	1288
1240	G.2 Graph-Aware Q&A		
1241	The Q&A system combines vector-based retrieval		
1242	with graph traversal. The EmbeddingStore in-		
1243	dexes both text chunks and node descriptions using		
1244	sentence-transformers (with a simple bag-of-words		
1245	fallback when unavailable). Given a question, the		
1246	GraphRetriever:		

1289 Each action is logged with timestamps, main-
1290 taining a complete edit history. Nodes carry
1291 a `verification_status` field (auto-generated,
1292 human-verified, human-edited, or flagged) that
1293 propagates through exports.

1294 **G.5 Export Formats**

1295 The system exports to multiple formats for different
1296 use cases:

- 1297 • **JSON:** Full graph data including nodes, edges,
1298 chunks, and metadata.
- 1299 • **Markdown:** Structured reading notes with
1300 section outlines.
- 1301 • **Mermaid:** Mind maps and flowcharts for vi-
1302 sualization.
- 1303 • **HTML:** Interactive D3.js-based graph visual-
1304 ization.

1305 All exports preserve traceability metadata, en-
1306 abling users to navigate from any extracted element
1307 back to the original source.

1308 **H Implementation and Deployment**

1309 The backend is implemented in Python with
1310 FastAPI for service endpoints and relies on stan-
1311 dard scientific libraries for retrieval and scoring
1312 (scikit-learn, NumPy, pandas). The multi-agent
1313 pipeline is defined in
1314 `textttbackend/agents/discovery/pca.py`, while the
1315 refactored deterministic pipeline is implemented in
1316 `textttbackend/core/paperfinder.py`. Both pipelines
1317 expose functionality through API servers, includ-
1318 ing a fast discovery variant designed for low-
1319 latency responses.

1320 The frontend is built with React and TypeScript
1321 and integrates discovery results through the API.
1322 Supabase provides authentication and persistent
1323 data storage for user profiles, communities, ses-
1324 sions, and paper metadata. Containerization sup-
1325 port is provided via a Dockerfile, and deployment
1326 configurations are included for common platforms
1327 (Railway, Render, and Vercel). Environment vari-
1328 ables control API URLs and database credentials,
1329 enabling local development or hosted deployment
1330 without code changes.