

# MOTIVGRAPH-SOIQ: INTEGRATING MOTIVATIONAL KNOWLEDGE GRAPHS AND SOCRATIC DIALOGUE FOR ENHANCED LLM IDEATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large Language Models (LLMs) hold substantial potential for accelerating academic ideation but face critical challenges in grounding ideas and mitigating confirmation bias for further refinement. We propose integrating motivational knowledge graphs and socratic dialogue to address these limitations in enhanced LLM ideation (MotivGraph-SoIQ). This novel framework provides essential grounding and practical idea improvement steps for LLM ideation by integrating a Motivational Knowledge Graph (MotivGraph) with a Q-Driven Socratic Ideator. The MotivGraph structurally stores three key node types—problem, challenge, and solution—to offer motivation grounding for the LLM ideation process. The Ideator is a dual-agent system utilizing Socratic questioning, which facilitates a rigorous refinement process that mitigates confirmation bias and improves idea quality across novelty, experimental rigor, and motivational rationality dimensions. On the ICLR25 paper topics dataset, MotivGraph-SoIQ exhibits clear advantages over existing state-of-the-art approaches across LLM-based scoring, ELO ranking, and human evaluation metrics.

## 1 INTRODUCTION

The potential of Large Language Models (LLMs) in supporting academic research (Achiam et al., 2023; Yang et al., 2024; Liu et al., 2024; Chen, 2024) within the domain of scholarly research has garnered increasing attention (Radensky et al., 2024; Si et al., 2024; Lu et al., 2024; Gupta & Pruthi, 2025). This includes the automated generation of literature reviews (Liang et al., 2025; Azaria et al., 2023), assistance in experimental design, and the enhancement of academic writing (Lu et al., 2024; Weng et al., 2024). Notably, leveraging the creativity of large language models to generate novel research ideas (Wang et al., 2024; Li et al., 2024; Si et al., 2024; Baek et al., 2024) is particularly compelling, which promises to accelerate the process of knowledge discovery, aiding researchers in transcending conventional thinking patterns and expanding the frontiers of exploration (Gottweis et al., 2025). However, the practical application of LLMs for generating research ideas still confronts two critical bottlenecks. Firstly, the generation process lacks a robust theoretical or factual grounding, which makes it challenging to create innovative and feasible ideas. Secondly, the issue of confirmation bias makes it difficult for LLMs to improve ideas.

**Motivation Grounding** Human researchers establish academic motivation connections through an extensive literature review, which helps uncover their underlying motivations and problem-solving approaches. This process enables them to navigate complex knowledge domains, understand fundamental concepts, and promote innovation across different disciplines. For example, researchers may observe that ant colonies utilize pheromone trails to identify optimal paths to food sources. Simultaneously, they recognize the challenge of optimizing data routing in large-scale wireless sensor networks. By linking these insights, they form an academic motivation connection between biological swarm intelligence and network optimization.

047 Such a connection can lead to the novel application of ant colony optimization algorithms to improve routing  
048 efficiency in sensor networks. The effectiveness of academic motivation connections lies in their ability to  
049 foster a comprehensive understanding of disparate fields and encourage combinatorial innovation, thereby  
050 generating novel and valuable ideas.

051 However, the internal knowledge of Large Language Models is probabilistic (Ye et al., 2025), unstable (Atil  
052 et al., 2024), and inherently biased due to training data distribution. Relying solely on large language models  
053 to generate “academic motivation connections” can lead to unreliable innovation. Concerns persist that  
054 LLM-generated ideas may be primarily hallucinatory, superficial (Gupta & Pruthi, 2025), or infeasible (Si  
055 et al., 2024). Although approaches have been proposed to ground LLMs with external academic resources  
056 for background information (Lu et al., 2024), their limited context windows hinder effective processing  
057 of extensive literature and the formation of deep connections. Consequently, enabling LLMs to generate  
058 innovative and feasible ideas necessitates a motivation knowledge base capable of providing a profound  
059 grasp of academic research’s underlying motivations and relationships, in a format compatible with LLM  
060 processing characteristics.

061  
062  
063 **Confirmation bias** Confirmation bias (Nickerson, 1998) is a cognitive bias where individuals favor information  
064 that confirms pre-existing beliefs (Wason, 1968). Human researchers are susceptible to favoring data that  
065 supports their hypotheses, sometimes overlooking contradictory evidence. Discussions between the mentor  
066 and the researcher are crucial for its mitigation in scientific contexts. In these settings, researchers present  
067 their hypotheses and reasoning to their mentors, who challenge assumptions, question methodologies, and  
068 highlight overlooked counterexamples, helping to correct biased reasoning and flawed assumptions. LLMs  
069 also exhibit this bias, struggling with novel thought generation and self-correction once an initial stance is  
070 established (Liang et al., 2023; Zhao et al., 2024). A key challenge in leveraging LLMs for academic ideation  
071 is enabling them to identify critical weaknesses in their generated ideas. While effective for superficial issues,  
072 current self-reflection methods fail to address fundamental shortcomings such as incorrect assumptions due to  
073 their vulnerability to confirmation bias (Liang et al., 2023). Thus, developing strategies for LLMs to refine  
074 ideas while actively mitigating this bias remains a considerable challenge.

075 In this paper, we propose Socratic LLM Ideation with Academic Motivation Graph (MotivGraph-SoIQ) to  
076 address the challenges above.

077  
078 **Contribution** Our main contributions are summarized as follows:

- 079 **1:** To address the lack of motivational grounding and limited self-improvement in LLM-based ideation, we  
080 propose **MotivGraph-SoIQ**. This unified framework integrates a Motivational Knowledge Graph with a  
081 Socratic ideation loop to produce grounded, high-quality ideas.
- 082  
083 **2:** We introduce **SciMotivMiner** to tackle the challenge of constructing a structured motivational resource  
084 from literature. SciMotivMiner automatically extracts (problem, challenge, method) triplets from published  
085 papers to build the MotivGraph, enabling motivational grounding for idea generation.
- 086  
087 **3:** We develop the **Q-Driven Socratic Ideator** to handle the difficulty of refining ideas and mitigating  
088 biases. This module employs a questioning-based self-improvement loop with four specialized tools  
089 for compelling graph exploration and strategic novelty injection, improving idea quality across multiple  
090 evaluation metrics.
- 091  
092 **4:** We conduct concise experiments on a topic set from ICLR25 papers, demonstrating that MotivGraph-SoIQ  
093 significantly outperforms strong baselines in novelty, experimental feasibility, motivational rationality, and  
diversity, achieving a 10.2 % improvement in novelty, a 6 % improvement in motivation, and an average  
ELO score 38 points higher.

## 2 METHOD

In this section, we detail our LLM-based ideation methodology, the MotivGraph-SoIQ Framework, which integrates two core components: (i) **MotivGraph**, a motivation-enhancing knowledge graph for structured motivation representation, and (ii) **Q-Driven Socratic Ideator**, an adversarial agentic system that refines ideas through “Socratic questioning” and “maieutics”.

### 2.1 MOTIVGRAPH

The **MotivGraph** serves two primary purposes. Firstly, it provides the underlying knowledge base to supply relevant knowledge crucial for the ideation process. Secondly, the explicit relationships between entities within the graph offer concrete examples of how problems can be framed and addressed. This structure is a valuable source of inspiration, specifically aiding LLMs in formulating clear and compelling motivations for novel research ideas.

#### 2.1.1 MOTIVGRAPH CONSTRUCTION

Amabile’s Componential Theory of Creativity (Amabile et al., 1996) posits that motivation constitutes one of the three essential components of innovation (alongside domain-relevant skills and creativity-relevant processes), with intrinsic motivation being particularly critical for breakthrough ideation. We design the **MotivGraph** as a graph structure consisting of three principal node types: *problem*, *challenge*, and *solution*. A *problem* node signifies a minimally granular research topic or task, a *challenge* node indicates a specific difficulty encountered within a *problem*, and a *solution* node represents a concrete method addressing a *challenge*. **Motivation** information is represented by triples formed through inter-node connections, specifically in the format (problem, challenge, solution). Each node is further characterised by two attributes: a concise and precise name for unique identification, and a description that provides further detail and aids in the graph’s semantic representation, matching, and retrieval processes.

The MotivGraph is represented as a graph  $G = (V, E)$ , where  $V$  is the set of nodes and  $E$  is the set of edges. The nodes  $V$  are classified into three types: problem ( $P$ ), challenge ( $C$ ), and solution ( $S$ ). The edge set  $E$  includes three distinct edge types: *parent-of* (for hierarchical links), *problem-challenge* (connecting  $P$  to  $C$ ), and *challenge-solution* (connecting  $C$  to  $S$ ). Figure 1 shows the construction process of MotivGraph. The specific construction method will be introduced in the following sections.

#### 2.1.2 SCIMOTIVMINER

For each scientific paper  $P$ , we employ our method, SciMotivMiner, denoted as  $SMM(P)$ , to process the paper and identify triples of related problems, challenges, and solutions. Consequently,  $SMM(P)$  outputs a set of  $n$  distinct (Problem, Challenge, Method) triples:  $\{(P_i, C_i, S_i)\}_{i=1}^n = SMM(P)$ .

For each extracted  $(P_i, C_i, S_i)$ , SciMotivMiner summarizes a concise entity name and a brief description. The naming process adheres to the rules to ensure clarity and consistency within the knowledge graph. The exact rules are detailed in the appendix B.3.

These stringent rules ensure a standardized, informative, and author-name-agnostic representation of research motivation and proposed solutions within the SciMotivMiner knowledge graph. For identical nodes, SciMotivMiner will merge them.

**Hierarchical Parent Node Addition** Academic problems and challenges inherently possess a hierarchical structure, with different papers addressing varying granularities. To capture these relationships and prevent knowledge fragmentation within the MotivGraph, we introduce Hierarchical Parent Node Addition for both

Problem (P) and Challenge (C) entities. This process organizes knowledge into a coherent hierarchy, crucial for practical exploration.

Our Parent Node Addition Algorithm operates iteratively. It begins by embedding all initial Problem/Challenge nodes into a vector space. The algorithm then repeatedly selects a focal node, identifies its k most similar neighbors within the current working set, and employs an LLM to evaluate their semantic coherence for merging. If the LLM deems an add appropriate, a new, more general parent node is created and linked to its children by parent-of edges. Processed nodes are then removed from the working set. This dynamic process ensures each node is considered for forming a parent at most once, building a multi-level abstract representation of the concepts. See the appendix B.4 for details.

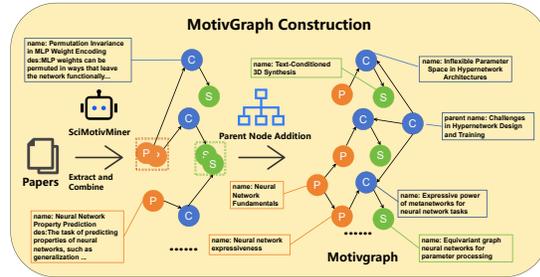


Figure 1: motivgraph construction pipeline

## 2.2 Q-DRIVEN SOCRATIC IDEATOR

The **Q-Driven Socratic Ideator** is a dual-agent system consisting of a mentor agent and a researcher agent. Its operational principles are inspired by Socratic questioning and maieutics. The mentor agent adheres to the elenchus (Socratic refutation) through triple-axis questioning—probing innovation, feasibility, and rationality—thereby exposing logical gaps without prescriptive solutions. The researcher agent operationalizes maieutics (intellectual midwifery) by synthesizing knowledge through: (1) introspective retrieval of dialogue history (“knowledge amniotic fluid”), and (2) external tool-augmented searches, ultimately ‘giving birth’ to refined ideas through self-directed epistemic labor. The following subsections detail the architecture, roles, and interaction dynamics of the agents within this system. The following sections delineate the two-phase architecture of the Q-Driven Socratic Ideator: (i) the **Exploration Phase**, and (ii) the **Deliberation Phase**.

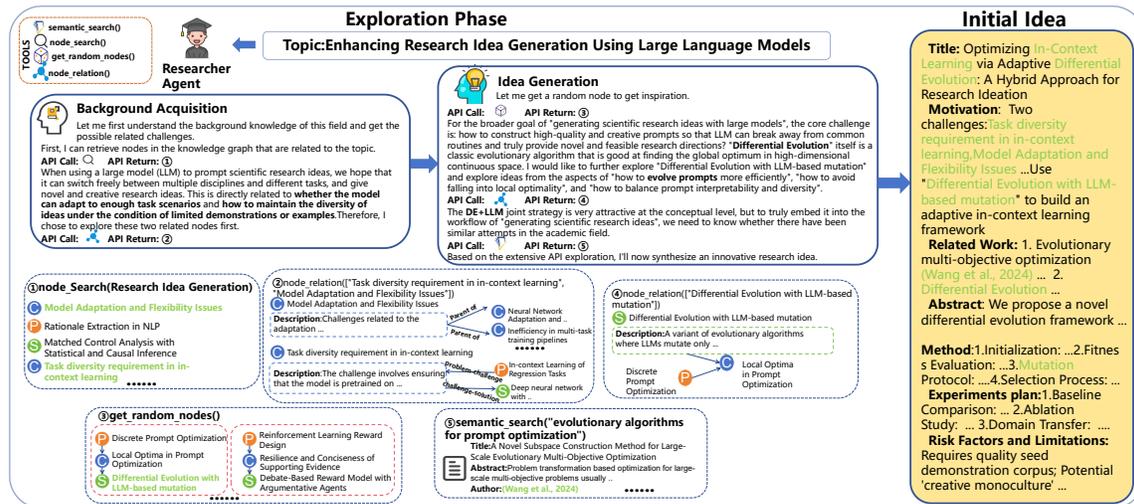


Figure 2: Exploration Phase Pipeline

### 2.2.1 EXPLORATION PHASE

The researcher agent primarily carries out the Exploration Phase. Based on the provided target domain or task description, the researcher agent performs knowledge exploration and generates innovative ideas. Figure 2 shows the process of the Exploration Phase. See Appendix B.6 for further details.

**(1) Knowledge Exploration and Ideation:** We designed three API tools to help the researcher agent better understand the target domain or task and generate an idea:

**(2) Graph Node Fuzzy Search:** This tool allows the researcher agent to obtain an overall understanding of the target domain/task by fuzzy-matching and retrieving related *problem*, *challenge*, and *solution* entities from the MotivGraph based on a search query.

**(3) Graph Node Relation Retrieval:** By providing an interesting node’s name, the researcher agent can retrieve its description and neighboring nodes, gaining hierarchical relationships and (problem, challenge, solution) motivation triplets. This deepens understanding and supports effective subsequent retrieval.

**(4) Semantic Scholar Literature Search:** This API provides query-based literature search, offering more specific information than the graph for a comprehensive understanding of particular challenges or technologies.

**(5) Get Random Nodes to Enhance Novelty:** After sufficient knowledge exploration, the researcher agent uses this API to obtain random problem-challenge-solution triples. It then attempts to apply these to the target domain, seeking potential connections or adaptations. This mechanism supports the "creativity-relevant processes" from Amabile’s Componential Theory of Creativity (Amabile et al., 1996), ensuring idea novelty. Simultaneously, the inherent logic of the MotivGraph’s (problem, challenge, solution) triples fosters "intrinsic motivation," driving the agent to explore adaptations of external nodes to the target domain, facilitating the discovery of new problems or innovative solutions.

### 2.2.2 DELIBERATION PHASE

Following the initial Exploration Phase, the researcher agent enters the Deliberation Phase, engaging in multi-round deliberation with the mentor agent. This phase rigorously evaluates and refines previously generated ideas through Socratic interaction (Figure 3).

During this phase, the mentor agent challenges the researcher agent’s idea from three angles—innovation, feasibility, and rationality—acting as a form of Socratic elenchus. It poses probing questions such as: “How does this solution transcend prior ideas?” (Innovation), “What tools would implement this?” (Feasibility), and “Why is your method effective?” (Rationality).

The researcher agent defends and justifies its idea using knowledge from the Exploration Phase. Through reflection, it may find flaws or gaps, prompting supplementary knowledge exploration via available API tools to address weaknesses. The refined idea and rationale are then presented for further questioning, forming a guided self-correction process—the essence of the maieutic method.

The mentor agent acts only as a critical questioner, facilitating refinement without offering answers or performing knowledge gathering. The number of deliberation rounds can be preset, and the mentor may end early if the idea is clearly strong or unviable. After deliberation, it gives a final evaluation: ACCEPT

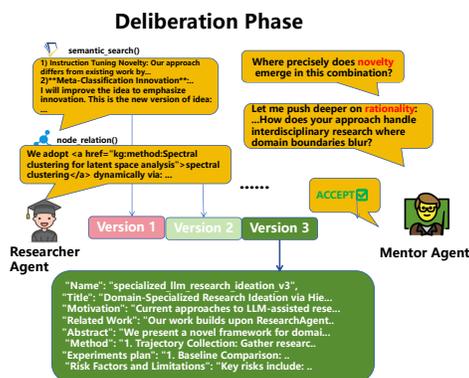


Figure 3: Deliberation Phase Pipeline

or REJECT. Rejected ideas are discarded to maintain quality and avoid retaining unjustified or incoherent concepts.

Baseline	Model	Diversity	LLM-evaluator			Human-evaluator			Length
			Nov.	Exp.	Moti.	Nov.	Exp.	Moti.	
AI-S-v2	DeepSeek-V3	0.27	7.22	8.07	8.21	6.35	6.25	5.85	2635
	Deepseek-R1	<b>0.52</b>	7.59	<b>8.36</b>	8.30	/	/	/	3013
	Qwen2.5-7B	0.24	6.10	7.22	7.28	/	/	/	3060
AI-Researcher	DeepSeek-V3	0.38	7.58	7.06	7.86	6.40	6.65	6.45	4985
	Deepseek-R1	0.32	7.94	7.65	8.19	/	/	/	4599
	Qwen2.5-7B	0.34	<b>7.14</b>	5.76	7.29	/	/	/	5465
SciPIP	DeepSeek-V3	0.42	7.61	7.23	7.61	6.20	/	5.05	4252
	Deepseek-R1	0.41	8.07	7.71	7.85	/	/	/	4230
	Qwen2.5-7B	0.35	6.51	6.04	6.46	/	/	/	5088
CycleResearcher	CycleResearcher-12B	0.29	6.61	7.52	7.39	5.50	6.25	5.35	7189
ResearchAgent	DeepSeek-V3	0.23	7.43	8.39	8.06	6.30	6.60	5.85	<b>15255</b>
	Deepseek-R1	0.25	8.02	8.33	8.17	/	/	/	<b>10204</b>
	Qwen2.5-7B	0.17	6.88	<b>7.67</b>	<b>7.60</b>	/	/	/	<b>13975</b>
Ours	DeepSeek-V3	<b>0.45</b>	<b>8.39</b>	<b>8.64</b>	<b>8.70</b>	<b>6.45</b>	<b>6.70</b>	<b>6.70</b>	4908
	deepseek-r1	0.45	<b>8.30</b>	8.00	<b>8.33</b>	/	/	/	4753
	qwen2.5-7b	<b>0.43</b>	6.46	6.64	6.52	/	/	/	3698
Real Paper	DeepSeek-V3	<b>1.00</b>	6.97	8.16	7.81	<b>7.08</b>	<b>7.36</b>	<b>8.05</b>	5030

Table 1: Evaluation Results: We use Fast-reviewer as LLM-evaluator. We manually evaluate and score ideas generated by DeepSeek-V3 using three dimensions: Novelty, experiment, and motivation. Ideas generated by SciPIP do not have experimental designs, so their experiments are not manually evaluated.

**Process Formalism** The iterative deliberation process and outcome can be formally represented. Let  $Idea_k$  be the state of the idea after round  $k$  ( $Idea_0$  is the initially generated idea), and  $I_k$  denote the interaction (mentor’s question and researcher’s response) at round  $k$ . The idea evolves based on the ideator agent’s function  $f_{\text{Researcher}}$ :

$$Idea_k = f_{\text{Researcher}}(Idea_{k-1}, I_k, \text{Exploration})$$

The deliberation phase concludes at round  $N_{\text{final}}$  (the maximum predefined rounds or an earlier termination round). The final evaluation,  $Eval$ , is given by the mentor agent’s function  $e_{\text{Mentor}}$  based on the final idea state and potentially the dialogue history:

$$Eval = e_{\text{Mentor}}(Idea_{N_{\text{final}}}, \text{Dialogue History})$$

where  $Eval \in \{\text{ACCEPT}, \text{REJECT}\}$ .

### 3 EXPERIMENT

To validate MotivGraph-SoIQ’s effectiveness, we conducted both comparative and ablation experiments. We constructed the MotivGraph and an evaluation dataset using publicly available literature. The MotivGraph provides motivational grounding, while the evaluation dataset, comprising 100 diverse ICLR 2025 paper topics and their core ideas, serves as ground truth for assessing generated ideas. Our comparisons against baselines demonstrate MotivGraph-SoIQ’s superiority, and ablation studies confirm the effectiveness of individual system components. See the Appendix B.5 for details on the dataset.

282 3.1 COMPARATIVE BASELINES  
283

284 To assess the effectiveness of our MotivGraph-SoIQ, we selected several baseline methods for comparison.  
285 The criteria for their selection were based on similarities in generating idea components similar to ours  
286 (Motivation, Related Work, Abstract, Method, Experiment Plan, Risk Factors, and Limitations) or employing  
287 entity/graph-based information enhancement. Please refer to Appendix B.1 for detailed information. The  
288 selected baselines are below:

289 **AI-Researcher:** This method, proposed in (Si et al., 2024), uses the author’s publicly available code to  
290 generate ideas.

291 **Cycle Researcher (12B):** Proposed in (Weng et al., 2024), we use the author’s publicly available code to  
292 generate idea proposals.

293 **AI-Scientist-v2:** This is an improved version of AI-Scientist (Yamada et al., 2025). We use the author’s  
294 publicly available code to generate ideas.

295 **SciPIP:** We use the author’s publicly available code to generate ideas.

296 **ResearchAgent:** We reproduce this method following the methodology described in the author’s paper (Baek  
297 et al., 2024).  
298

299 3.2 ABLATION STUDIES  
300

301 We conducted a series of ablation studies to better understand each component’s contribution to our method’s  
302 overall performance and validate the necessity of these design choices. This section systematically removed  
303 or modified one or more parts of the MotivGraph and the Critique-Driven Agent System. We tested these  
304 variants using the same experimental setup and evaluation metrics as the whole method. By comparing the  
305 performance of different variants, we can quantify the effectiveness of each component and reveal the key  
306 roles they play in the ideation process.

307 The following ablation variants were tested:

308 1. **W/O Mentor:** The deliberation loop involving the mentor agent was removed in this configuration. The  
309 researcher agent generates and revises the idea by themselves.

310 2. **W/O Graph:** In this experimental condition, we intercept all MotivGraph API calls and return complete  
311 texts or abstracts from the corpus of research papers used to build the knowledge graph. The researcher  
312 agent’s access to knowledge is thus limited to this simulated interface, which provides document-level outputs  
313 rather than structured graph relationships.

314 3. **SCI-PIP Graph W/O Mentor:** This variant replaced our MotivGraph with the “concept-paper” graph  
315 constructed in SciPIP (Wang et al., 2024). SciPIP’s built-in retriever was used to retrieve relevant entities  
316 from its graph structure, which were then used for knowledge augmentation.

317 4. **W/O Graph + W/O Mentor:** In this variant, neither the MotivGraph nor the mentor agent’s deliberation  
318 process was utilized.

319 5. **W/O Semantic Scholar:** This variant retained the Semantic Scholar API for metadata retrieval but  
320 constrained its output to paper titles only, rather than complete metadata(including title, abstract, author, and  
321 publication year).  
322

323 3.3 EVALUATION SETUP  
324

325 Given the time-consuming and subjective nature of manual evaluation, and the documented efficacy of LLMs  
326 in judging text quality (Zheng et al., 2023; Fu et al., 2023; Liu et al., 2023), we adopted a model-based  
327  
328

evaluation approach. This includes LLM direct evaluation and Swiss Tournament evaluation. For diversity assessment, we calculate diversity as **1-MeanSimilarity** among multiple ideas generated for the same topic (Si et al., 2024). See the Appendix B.7 for details.

		Nov.	Moti.	Exp.	Average
Model	Ours	<b>1072</b>	<b>1061</b>	<b>1061</b>	<b>1064</b>
	AI-Scientist-v2	1034	1016	1028	1026
	ResearchAgent	1002	1011	1002	1005
	AI-Researcher	1012	995	1001	1003
	RealPaper	980	1020	1004	1001
	SciPIP	1018	982	1002	1000
	CycleResearcher	879	912	899	897
Human	Ours	1038	1024	1026	1029
	RealPaper	<b>1071</b>	<b>1064</b>	<b>1063</b>	<b>1066</b>
	AI-Researcher	1013	1015	1020	1016
	AI-Scientist-v2	1010	1005	1013	1009
	ResearchAgent	990	1003	1012	1002
	SciPIP	1008	987	977	991
	CycleResearcher	966	988	983	979

Table 2: Comparison of Ideation Methods

### 3.4 IMPLEMENT

We selected three models—Qwen2.5-7B-Instruct (Qwen et al., 2025), DeepSeek-V3, and DeepSeek-R1 (Guo et al., 2025)—to investigate how models with different capabilities affect idea generation methods. Using a dataset of topics extracted from papers accepted at ICLR 2025, we generated at least three ideas per topic with each technique. Subsequently, we calculated the diversity of the generated ideas and employed Fast-Reviewer to quickly evaluate these ideas based on three dimensions: Novelty (Nov.), Experiment (Exp.), and Motivation (Moti.).

Additionally, we use DeepSeek-V3 to conduct a Swiss Tournament evaluation on the generated ideas across the Novelty, Motivation, and Experiment dimensions, computing ELO scores for each dimension and an overall average score.

To further ensure the reliability of our evaluation, we replaced the automated Swiss Tournament assessment and LLM assessment with manual evaluations and reported corresponding ELO scores with direct scores. Since manual evaluation is time-consuming and labour-intensive, we only selected ideas generated by DeepSeek-V3, chosen topics, and selected one idea per topic for manual evaluation.

## 4 RESULT AND ANALYSIS

### 4.1 COMPARATIVE BASELINES

Table 1 presents the comparative results with the baselines. Experimental results show that our method has obvious advantages when DeepSeek-V3 generates ideas. Regarding diversity, Novelty, Experiment, and Motivation, our process is 0.03, 0.78, 0.25, and 0.49, higher than the second-best baseline regarding automatic evaluation. Manual evaluation results show that our method is 0.05, 0.05, and 0.25 higher than the second-best baseline regarding Novelty, Experiment, and Motivation. When the Qwen2.5-7B small parameter model is used, the model’s ability to call APIs and integrate API return information is insufficient, and the number of API calls is abnormally high or low. At the same time, the context length that the small model can use is inadequate. In multiple rounds of modifications, part of the historical records often need to be discarded,

376 which reduces the quality of idea generation to a certain extent. As for DeepSeek-R1, we can see that the  
 377 Novelty and Motivation scores of the idea are still high due to the existence of the graph, but the scores of the  
 378 three dimensions are lower than those of DeepSeek-V3. This is because the reasoning model requires long  
 379 thinking, so the API call is planned before the API returns the result, which hinders the model from gradually  
 380 exploring in depth.

381 Table 2 compares the ELO score with the baseline. The results show that our method scores 28 points, 45  
 382 points, and 33 points higher than the second-best method (except Real Paper) in novelty, motivation, and  
 383 experiments, respectively, and an average score of 38 points higher. The ELO scores of human evaluation are  
 384 25 points, 9 points, and 6 points higher in Novelty, Experiment, and Motivation, respectively.

## 386 4.2 ABLATION STUDIES

387 We conducted ablation studies to evaluate the contribution of key components in our framework. Table 3  
 388 summarizes the results.

389 First, we examined the role of the knowledge graph. Removing its hierarchical structure and using only raw  
 390 text reduced the Novelty and Experiment scores by 0.4 and 0.2, respectively. Replacing it with a generic  
 391 scipip-graph baseline also degraded performance (-0.3 Novelty, -0.2 Experiment, -0.1 Motivation), confirming  
 392 the effectiveness of our customized graph design.

393 Next, removing the mentor interaction phase caused larger drops across all metrics (-0.9 Novelty, -0.8  
 394 Experiment, -0.6 Motivation), highlighting the importance of iterative discussion and refinement with the  
 395 mentor.

396 Interestingly, disabling both the mentor and graph (w/o mentor + w/o graph) yielded slightly higher scores  
 397 than the w/o mentor setup alone. We hypothesize this occurs because, without the graph’s divergent influence,  
 398 the model generates more conventional yet stable ideas. Figure 4 illustrates the final score versus discussion  
 399 rounds.

400 Finally, removing Semantic Scholar content except titles reduced Novelty and Experiment scores, showing  
 401 that detailed background knowledge enhances innovation and experimental soundness.

## 405 5 CONCLUSION

406 LLMs offer great promise for academic ideation but  
 407 face challenges with idea grounding and confirma-  
 408 tion bias. We introduce MotivGraph-SoIQ, a novel  
 409 framework that enhances LLM ideation by integrat-  
 410 ing a Motivational Knowledge Graph for ground-  
 411 ing from literature and a Q-Driven Socratic Ideator.  
 412 This dual-agent system uses Socratic questioning to  
 413 refine ideas, mitigating confirmation bias and im-  
 414 proving novelty, experimental feasibility, and moti-  
 415 vation. Our results demonstrate MotivGraph-SoIQ’s  
 416 effectiveness and superior performance across LLM-  
 417 based scoring, ELO ranking, and human evaluation.  
 418 Ablation studies confirm the crucial contributions of  
 419 both MotivGraph and the Socratic dialogue. This  
 420 work highlights the power of combining structured  
 421 knowledge with interactive, critique-based refine-  
 422 ment for robust LLM ideation.

Methods	Nov.	Exp.	Moti.
Ours	<b>8.4/6.5</b>	<b>8.6/6.7</b>	<b>8.7/6.7</b>
- w/o graph	8.0/5.7	8.40/6.2	8.7/5.7
- w/ scipip-graph	8.1/5.8	8.4/6.2	8.6/6.1
- w/o mentor	7.5/5.7	7.8/5.5	8.1/5.7
- w/o mentor & graph	7.7/5.7	8.2/5.7	8.5/5.9
- w/o semantic scholar	8.1/6.0	8.4/6.4	8.7/6.3

Table 3: Results of ablation study on references and entities. The scores on the left of “/” are obtained using Fast-Reviewer evaluation, and those on the right are obtained by manual evaluation.

## 6 LIMITATIONS

While our findings are promising, we acknowledge several limitations in the current work. The scope of our constructed MotivGraph is presently limited, primarily encompassing knowledge within the AI domain and lacking comprehensive coverage of other scientific disciplines. Expanding its domain coverage is essential for realising the full potential of cross-disciplinary idea generation. However, our constructed MotivGraph holds considerable potential for uncovering connections across diverse scientific disciplines and presenting these associations to large language models for their utilisation. Furthermore, due to constraints on available resources and time, our experimental validation was conducted on a specific dataset size, and we evaluated the framework using a limited variety of LLM models. Future work should focus on scaling up the experimental evaluation to a larger dataset and testing a more diverse range of underlying LLMs to confirm the generalizability of our findings.

For future research, we also plan to explore extending the MotivGraph to incorporate other academic knowledge and relationships. Further investigation into alternative dialogue strategies within the Socratic framework could yield additional insights.

## 7 ETHICS STATEMENT

Our system is developed with the explicit and sole purpose of serving as an assistive tool to augment human creativity and facilitate the discovery of novel research ideas within the academic domain. Our goal is to empower researchers by providing inspiration, helping to overcome ideation blocks, and suggesting potentially fruitful avenues for investigation grounded in existing knowledge.

We unequivocally condemn and strongly disavow any potential misuse of this system. This includes, but is not limited to, using the system to generate ideas or methods for illegal activities, unethical research practices, harmful technologies, malicious applications, or any purpose that could cause societal harm, violate privacy, or infringe upon human rights. Users are solely responsible for the evaluation, validation, and ethical implications of any system-generated idea and its subsequent application. The system is designed to be a creative aid, not an autonomous decision-maker or a substitute for human ethical reasoning and responsibility.

## 8 ACKNOWLEDGEMENT

This work is supported by the National Natural Science Foundation of China (No.U24A20335, No. 62176257, No.62576340). This work is sponsored by Beijing Nova Program (No.20250484750) and supported by Beijing Natural Science Foundation (L243006). This work is also supported by the Youth Innovation Promotion Association CAS.

## REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Teresa M Amabile et al. *Creativity and innovation in organizations*, volume 5. Harvard Business School Boston, 1996.
- Berk Atil, Alexa Chittams, Liseng Fu, Ferhan Ture, Lixinyu Xu, and Breck Baldwin. Llm stability: A detailed analysis with some surprises. *arXiv preprint arXiv:2408.04667*, 2024.

- 470 Amos Azaria, Rina Azoulay, and Shulamit Reches. Chatgpt is a remarkable tool – for experts, 2023. URL  
471 <https://arxiv.org/abs/2306.03102>.
- 472  
473 Jinheon Baek, Sujay Kumar Jauhar, Silviu Cucerzan, and Sung Ju Hwang. Researchagent: Iterative research  
474 idea generation over scientific literature with large language models. *arXiv preprint arXiv:2404.07738*,  
475 2024.
- 476 Huajun Chen. Large knowledge model: Perspectives and challenges. *Data Intelligence*, 6(3):587–620, August  
477 2024. ISSN 2096-7004. doi: 10.3724/2096-7004.di.2024.0001. URL [http://dx.doi.org/10.3724/  
478 2096-7004.di.2024.0001](http://dx.doi.org/10.3724/2096-7004.di.2024.0001).
- 479 Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. Gptscore: Evaluate as you desire. *arXiv preprint  
480 arXiv:2302.04166*, 2023.
- 481  
482 Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom Myaskovsky,  
483 Felix Weissenberger, Keran Rong, Ryutaro Tanno, et al. Towards an ai co-scientist. *arXiv preprint  
484 arXiv:2502.18864*, 2025.
- 485  
486 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma,  
487 Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement  
488 learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 489  
490 Tarun Gupta and Danish Pruthi. All that glitters is not novel: Plagiarism in ai generated research. *arXiv  
491 preprint arXiv:2502.16487*, 2025.
- 492  
493 Long Li, Weiwen Xu, Jiayan Guo, Ruochen Zhao, Xingxuan Li, Yuqian Yuan, Boqiang Zhang, Yuming Jiang,  
494 Yifei Xin, Ronghao Dang, et al. Chain of ideas: Revolutionizing research via novel idea development with  
495 llm agents. *arXiv preprint arXiv:2410.13185*, 2024.
- 496  
497 Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and  
498 Zhaopeng Tu. Encouraging divergent thinking in large language models through multi-agent debate. *arXiv  
499 preprint arXiv:2305.19118*, 2023.
- 500  
501 Xun Liang, Jiawei Yang, Yezhaohui Wang, Chen Tang, Zifan Zheng, Shichao Song, Zehao Lin, Yebin Yang,  
502 Simin Niu, Hanyu Wang, et al. Surveyx: Academic survey automation via large language models. *arXiv  
503 preprint arXiv:2502.14776*, 2025.
- 504  
505 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng,  
506 Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- 507  
508 Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-eval: Nlg evaluation  
509 using gpt-4 with better human alignment. *arXiv preprint arXiv:2303.16634*, 2023.
- 510  
511 Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist: Towards  
512 fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*, 2024.
- 513  
514 Raymond S Nickerson. Confirmation bias: A ubiquitous phenomenon in many guises. *Review of general  
515 psychology*, 2(2):175–220, 1998.
- 516  
517 Qwen, ., An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,  
518 Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang,  
519 Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li,  
520 Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang  
521 Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang,  
522 and Zihan Qiu. Qwen2.5 technical report, 2025. URL <https://arxiv.org/abs/2412.15115>.

- 517 Marissa Radensky, Simra Shahid, Raymond Fok, Pao Siangliulue, Tom Hope, and Daniel S Weld. Scideator:  
518 Human-llm scientific idea generation grounded in research-paper facet recombination. *arXiv preprint*  
519 *arXiv:2409.14634*, 2024.
- 520 Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv*  
521 *preprint arXiv:1908.10084*, 2019.
- 522 Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. Can llms generate novel research ideas? a large-scale  
523 human study with 100+ nlp researchers. *arXiv preprint arXiv:2409.04109*, 2024.
- 524 Wenxiao Wang, Lihui Gu, Liye Zhang, Yunxiang Luo, Yi Dai, Chen Shen, Liang Xie, Binbin Lin, Xiaofei  
525 He, and Jieping Ye. Scipip: An llm-based scientific paper idea proposer. *arXiv preprint arXiv:2410.23166*,  
526 2024.
- 527 Peter C Wason. Reasoning about a rule. *Quarterly journal of experimental psychology*, 20(3):273–281, 1968.
- 528 Yixuan Weng, Minjun Zhu, Guangsheng Bao, Hongbo Zhang, Jindong Wang, Yue Zhang, and Linyi Yang.  
529 Cycleresearcher: Improving automated research via automated review. *arXiv preprint arXiv:2411.00816*,  
530 2024.
- 531 Yutaro Yamada, Robert Tjarko Lange, Cong Lu, Shengran Hu, Chris Lu, Jakob Foerster, Jeff Clune, and  
532 David Ha. The ai scientist-v2: Workshop-level automated scientific discovery via agentic tree search. *arXiv*  
533 *preprint arXiv:2504.08066*, 2025.
- 534 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu,  
535 Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- 536 Xiaotian Ye, Mengqi Zhang, and Shu Wu. Open problems and a hypothetical path forward in llm knowledge  
537 paradigms. *arXiv preprint arXiv:2504.06823*, 2025.
- 538 Suifeng Zhao, Tong Zhou, Zhuoran Jin, Hongbang Yuan, Yubo Chen, Kang Liu, and Sujian Li. Awecita:  
539 Generating answer with appropriate and well-grained citations using llms. *Data Intelligence*, 6(4):1134–  
540 1157, 2024.
- 541 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,  
542 Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena.  
543 *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.

## 544 A FURTHER ANALYSIS

545 This subsection discusses the intermediate results produced during the idea generation process.

546 **Idea score vs. the number of rounds.** Figure 4 illustrates the relationship between the final idea score  
547 and the number of discussion rounds. From this figure, it can be observed that discussion contributes to an  
548 improvement in overall quality. Furthermore, a higher initial quality often correlates with fewer discussion  
549 rounds, and scores are notably higher when the mentor raises fewer questions. Nevertheless, engaging in  
550 more discussion rounds can also enhance the overall quality of the ideas.

551 **API Usage.** Figures 4 present the frequency and distribution of API calls made by the researcher agent  
552 across different rounds during the idea generation process, respectively. These figures demonstrate that our  
553 constructed researcher agent can autonomously invoke tools and independently determine tool usage based  
554 on the specific problem context.

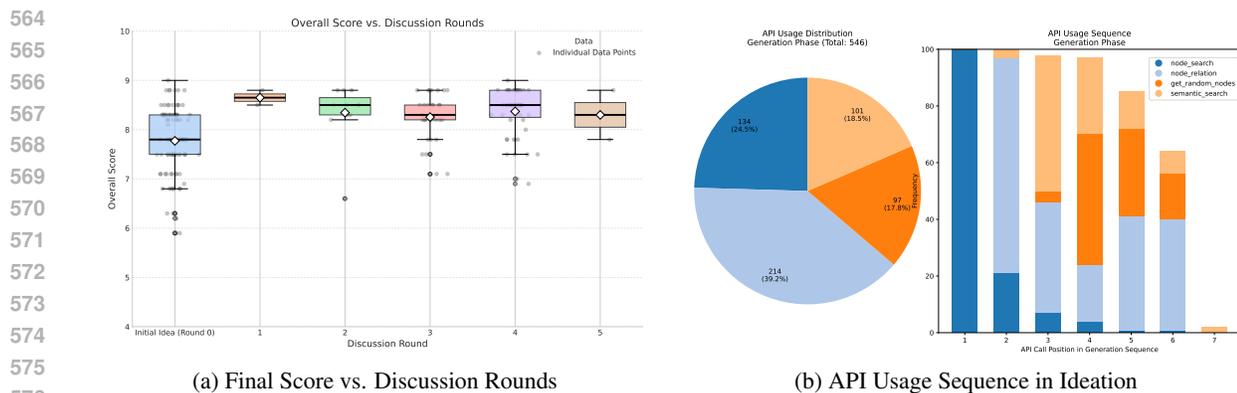


Figure 4: Combined Analysis of Discussion Rounds and API Usage Patterns

**Differences between backbone models.** As shown in the table 4, for the Qwen model, increasing the parameter count from 7 B to 32 B yields a marked improvement in idea quality. Overall, our method’s performance can benefit from stronger model capabilities; however, when the model size reaches 72 B, quality actually declines. Our observations reveal that Qwen-2.5-72 B begins to produce garbled output under long-context conditions, which we believe indicates a sharp drop in its comprehension and reasoning capabilities once the input exceeds a certain length. We observed the same behavior with Qwen3. Indeed, Qwen3 generated extensive mixed-language garble that prevented the pipeline from functioning correctly. Consequently, we conclude that Qwen models show a substantial performance gap compared to DeepSeek when tasked with understanding and analyzing large volumes of text.

Model	Novelty	Motivation	Experiment
DeepSeek-V3	8.39	8.70	8.64
Qwen-7B	6.46	6.52	6.44
Qwen-14B	6.55	6.95	7.33
Qwen-32B	6.85	7.95	7.70
Qwen-72B	6.90	7.05	6.55

Table 4: LLM-evaluator scores for different Qwen model sizes and DeepSeek-V3.

**Generalizability to Other Scientific Domains.** Theoretically, our method, MotivGraph-SoIQ, offers strong generalizability across disciplines. Its MotivGraph component supplies the large model with a motivational foundation for idea conception in <Problem, Challenge, Solution>, reflecting a basic scientific-research paradigm in many fields. Moreover, using Socratic dialogue to refine ideas iteratively is likewise a common research practice.

We collected 185 recent papers from high-quality medical journals (Nature Medicine, Nature Biomedical Engineering, and IEEE Transactions on Medical Imaging) to validate our approach in another domain empirically. We clustered these into 30 topics for idea generation and used the remaining 155 papers to construct a small-scale knowledge graph. We then compared our method against two strong baselines:

ResearchAgent and CycleResearcher (a domain-knowledge fine-tuned model). Because our original evaluator was trained only on AI-domain papers, we replaced it here with DeepSeek-r1, which offers comparable performance. Table 5 shows that MotivGraph-SoIQ continues to perform effectively in the medical domain.

Model	Novelty	Motivation	Experiment
Ours	8.04	8.72	8.02
ResearchAgent	7.55	8.27	8.01
CycleResearcher	6.85	7.87	6.89

Table 5: LLM-evaluator scores for different methods in medical domain.

## B DETAILS

### B.1 BASELINE IMPLEMENT

**AI-Researcher:** This method represents a simple yet effective LLM ideation approach that integrates Retrieval-Augmented Generation (RAG), filters duplicate ideas using vector similarity, and employs an LLM-based automatic ranker inspired by a Swiss-tournament design. It is a typical example of using a single LLM agent for idea generation without explicitly constructing a knowledge graph. Comparing with AI-Researcher helps us understand the performance level of a general or relatively straightforward LLM generation agent in the context of ideation. The ideas it generates primarily include “Motivation”, “Proposed Method”, “Step-by-Step Experiment Plan,” and “Test Case Examples,” which share similarities in structure with the ideas produced by our method.

**Cycle Researcher:** This method introduces Iterative Preference Training, leveraging extensive prior literature and review feedback to train the Cycle Researcher model. It can generate paper proposals covering motivation, idea (method), and experimental setup, a structure akin to our generated ideas. We use Cycle Researcher as a baseline to assess the ideation capability of LLMs trained through reinforcement learning. In our experiments, we opted for the 12B model for comparison primarily to conserve idea generation time and resources. As indicated in the original Cycle Researcher paper, the 12B model exhibited performance comparable to, and in many metrics even superior to, their 123B model for this task. Additionally, we replaced their original Bib literature database with the Semantic Scholar API for the RAG component.

This is a section in the appendix.

**AI-Scientist-v2:** This is an improved version of AI-Scientist (Yamada et al., 2025), which enriches the content of generated ideas and integrates Semantic Scholar as a function call. Similar to our method’s approach to external information retrieval, AI-Scientist-v2 utilizes external knowledge via API calls. In generating ideas for comparison using AI-Scientist-v2, we set the number of reflection steps to 5 and generated five ideas per topic.

**SciPIP:** The core methodology of SciPIP (Wang et al., 2024) lies in combining multi-angle literature retrieval and a dual-path idea generation strategy. It first retrieves literature content and related entities based on the provided topic and then generates ideas through brainstorming and RAG. Similar to our method, SciPIP constructs a knowledge graph, specifically a “concept-paper” graph where concepts are extracted from papers

658 by a large model. However, it does not explicitly structure knowledge around challenges and solutions. This  
 659 method serves as an excellent comparison point to demonstrate the effectiveness of our Challenge-Solution  
 660 Knowledge Graph. We used SciPIP for standalone idea generation and integrated its graph entity retrieval  
 661 module to replace our Challenge-Solution Knowledge Graph during the idea generation process to compare  
 662 the graph structures directly. We use the dual-path approach mentioned in the paper for idea generation.  
 663

664 **ResearchAgent:** ResearchAgent (Baek et al., 2024) is a system designed to assist researchers in iterative  
 665 research idea generation using Large Language Models (LLMs). It aims to produce novel and impactful  
 666 research ideas by augmenting information from scientific literature and employing collaborative LLM-driven  
 667 review agents for iterative optimization. Its strategy of information enhancement via an “academic graph +  
 668 entity knowledge storage” and iterative optimization through a multi-agent collaborative review loop shares  
 669 similarities with our method but presents distinct differences, making it a valuable subject for comparative  
 670 experiments.  
 671

## 672 B.2 FAST-REVIEWER TEST:

674 We tested Fast-Reviewer Test on our constructed dataset, measuring AUC scores for positive-negative  
 675 discrimination across Novelty, Experiment, and Motivation categories. Table 6 shows the AUC scores.  
 676

## 677 B.3 SCIMOTIVMINER RULES:

679 **Problem Entity ( $P_i$ ):** The name of the **Problem**  
 680 entity ( $P_i$ ) represents the overall research task, do-  
 681 main, or high-level objective that the paper ad-  
 682 dresses. The naming follows the structure ‘[Gen-  
 683 eral Task/Field of Study]’ and aims for **3-7 words**.  
 684 These names must be generalized and **strictly avoid**  
 685 authors’ specific, non-generalized names or abbrev-  
 686 viations. Problem entity names are derived solely  
 687 from the context in the paper’s Introduction section.  
 688

689 The corresponding problem description provides a  
 690 brief (ideally **1-2 sentences**), neutral, high-level def-  
 691 inition of the overall research task, field of study,  
 692 or objective. This description focuses solely on the  
 693 task or field or its general purpose/goals, presented as a standalone concept. Crucially, this description **must**  
 694 **not** include any mention of challenges, limitations, difficulties, or specific areas of focus motivated by these  
 695 challenges. It is formulated as a universal definition, informed by the Introduction section.  
 696

697 **Challenge Entity ( $C_i$ ):** The name of the **Challenge** entity ( $C_i$ ) captures a specific, atomic difficulty,  
 698 limitation, gap, or existing shortcoming *within the identified Problem* that the paper aims to address. Its  
 699 naming strictly adheres to the structure ‘[Specific Difficulty/Limitation] in [Aspect of Problem/Domain  
 700 Context]’ to clearly state the precise difficulty and its context within the problem or domain. These names aim  
 701 for **5-8 words**, prioritizing the required structure and specificity. As with problem entities, authors’ specific,  
 702 non-generalized names or abbreviations for challenges are **strictly avoided**, and names are derived from the  
 703 Introduction section.  
 704

705 The challenge description, summarized from the Introduction section (ideally **2-3 sentences**), explains this  
 706 specific difficulty, limitation, or gap and details how it relates to the broader Problem.  
 707

Model	AUC		
	Novelty	Motivation	Experiment
Fast-Reviewer	<b>0.76</b>	0.56	0.66
DeepSeek-V3	0.68	0.57	0.53
deepseek-R1	0.75	0.58	<b>0.70</b>

Table 6: AUC Score

705 **Solution Entity ( $S_i$ ):** For the **Solution** entity ( $S_i$ ), the name captures the essential technical approach,  
 706 category, or fundamental principle employed to address the Challenge. A crucial constraint is that the authors'  
 707 specific name, acronym, or code name for their proposed solution (or any non-generalized term they introduce)  
 708 is **strictly not used** in the entity name, drawing instead on general technical terms or descriptions of the  
 709 solution's core components or principles. Solution names aim for **7-10 words**. For solutions described with a  
 710 citation in the Introduction, their established general name or common abbreviation (if widely recognized  
 711 and within the word count aim) is used, based on the Introduction description. For novel solutions (typically  
 712 described without a citation in the Introduction), the solution section is consulted to understand the core  
 713 technical approach and fundamental principles, and the name is generated using general technical terms or  
 714 essential component descriptions based on this technical understanding from both sections.

715 The solution description provides a brief (ideally **2-3 sentences**) explanation of the solution's core technical  
 716 aspects, focusing on *how* it works technically. If the Introduction's description is high-level, results-focused,  
 717 or lacks sufficient technical detail, the solution section is consulted to incorporate key technical aspects  
 718 explaining the approach.

#### 721 B.4 DETAILED HIERARCHICAL PARENT NODE ADDITION

722  
 723 Following the extraction process from the papers, we obtain a set of  $n$  distinct knowledge triplets,  
 724  $(P_i, C_i, S_i)_{i=1}^n$ . While initially extracted as independent triples, the problems and challenges described  
 725 within them exhibit inherent relationships. Academic problems inherently possess a hierarchical structure, and  
 726 different papers address problems at varying granularities. For example, one paper might focus on the broad  
 727 area of 'Machine Translation', while another delves into 'Low-Resource Machine Translation for Indigenous  
 728 Languages'. To capture these relationships and further associate the knowledge, we construct a hierarchical  
 729 structure for the graph by introducing parent nodes for both Problem ( $P$ ) and Challenge ( $C$ ) entities. This  
 730 hierarchical organisation is crucial to prevent the knowledge base from becoming overly fragmented or  
 731 unstructured, making it challenging to comprehend and navigate. Without this hierarchy, the graph would  
 732 fail to fully leverage its advantages for thoroughly organizing complex information, hindering compelling  
 733 exploration during subsequent ideation processes.

734 To acquire these parent nodes and establish hierarchical relationships within the sets of Problem ( $P$ ) and  
 735 Challenge ( $C$ ) nodes, we propose the **Parent Node Addition Algorithm**. This process is applied separately  
 736 to the Problem ( $P$ ) and Challenge ( $C$ ) node collections.

737 All original Problem and Challenge nodes are initially embedded into a vector space to enable subsequent  
 738 similarity search based on their semantic representations. This vector space representation is fundamental for  
 739 quantifying the semantic relationships between nodes.

740 The algorithm operates on an initial set  $S$ , which at the start of the process, contains all nodes from either the  
 741 Problem or Challenge set being processed. The core mechanism involves iteratively processing nodes within  
 742 this set  $S$  until it becomes empty.

743 The algorithm maintains  $S$  as a dynamic working set. It repeatedly selects a node  $N$  from the current set  
 744  $S$ . For this focal node  $N$ , the algorithm identifies its  $k$  most similar neighbours based on the pre-calculated  
 745 vector embeddings. A critical filtering step is then applied: only those similar neighbors that *also* remain  
 746 present in the current working set  $S$  are retained as potential candidates for grouping with  $N$ . Let this filtered  
 747 set of eligible similar nodes be  $V_{filtered}$ .

748 A Large Language Model (LLM) is crucial at this stage. It evaluates the semantic coherence and potential for  
 749 forming a higher-level concept when considering the focal node  $N$  and the nodes in  $V_{filtered}$ . Based on this  
 750 evaluation, the LLM decides whether a merge operation should occur.  
 751

752 A new parent node is created if the LLM determines that a merge is appropriate and the set  $V_{filtered}$  is not  
753 empty. This new node represents a more general theme or domain that encapsulates the concepts expressed  
754 by  $N$  and the nodes in  $V_{filtered}$ . Directed edges, labelled *parent-of*, are added from this new parent node to  
755  $N$  and to every node  $v \in V_{filtered}$ , establishing their hierarchical link.

756 Following the decision and potential merge, the current node  $N$  is removed from the set  $S$ , as it has been  
757 processed in this iteration. Furthermore, if a merge occurred, all the nodes in  $V_{filtered}$  that became children  
758 of the new parent node are also removed from the set  $S$ . This dynamic update ensures that each node is  
759 considered for forming a parent at most once in this pass and that nodes already integrated into a higher level  
760 via merging are no longer candidates within the same pass.

761 The iterative selection and processing of nodes from the set  $S$  continues until  $S$  becomes empty. At this  
762 point, all nodes from the initial set have been either processed as a focal node or removed because they were  
763 merged as children. The parent nodes created during this process represent a higher level of abstraction for  
764 the grouped concepts within the original set  $S$ .

## 765 B.5 DATASET CONSTRUCTION AND EVALUATION DETAILS

766 MotivGraph Dataset We constructed the MotivGraph from 8625 accepted papers from ICLR 2024, ICML  
767 2024, and NeurIPS 2024, collected from OpenReview and other sources. Using the SciMotivMiner method  
768 (detailed in Section 2.1.1) with DeepSeek-V3 as the extractor on the full text of these papers, we obtained  
769 25515 solution nodes, 31158 challenge nodes, and 12137 problem nodes. Node descriptions were vectorized  
770 using all-MiniLM-L6-v2 (Reimers & Gurevych, 2019). Subsequently, the Hierarchical Parent Node Addition  
771 method (detailed in Section B.4) established 37367 PARENT\_OF relationships and added 7089 parent nodes.

772 Evaluation Dataset We clustered the titles of all accepted ICLR 2025 papers for the evaluation dataset using  
773 all-MiniLM-L6-v2. From these clusters, we selected 100 papers(excluding papers used for Fast-Reviewer  
774 training) representing diverse topics. DeepSeek-V3 (Liu et al., 2024) extracted each selected paper’s core  
775 idea and topic. The extracted core ideas served as ground truth, matching our method’s output format for  
776 subsequent comparisons, while the extracted topics served as input for the idea generation process.

## 777 B.6 DETAILED RESEARCHER AGENT’S TOOLSET

778 The researcher agent within our Q-Driven Socratic ideator has four specialized tools to facilitate comprehen-  
779 sive knowledge exploration and foster innovative ideation.

### 780 GRAPH NODE FUZZY SEARCH

781 The researcher agent provides a search query. This API returns the names and types of the top K similar nodes  
782 based on the semantic similarity between the search query and the descriptions of nodes in the Motivational  
783 Knowledge Graph (MotivGraph). This tool enables the researcher Agent to gain an overarching understanding  
784 of the target domain or task by identifying related *problem*, *challenge*, and *solution* entities within that  
785 domain.

### 786 GRAPH NODE RELATION RETRIEVAL

787 Given the name of an interesting node, this API returns the node’s description, names, and types of its  
788 neighboring nodes. The researcher agent can retrieve hierarchical relationships between nodes and the  
789 (problem, challenge, solution) triplets representing critical motivational information through this tool. This  
790 contextual information deepens the researcher agent’s understanding of the target domain/task, facilitates  
791 more effective subsequent retrieval, and establishes a robust foundation for the ideation phase.

## 799 SEMANTIC SCHOLAR LITERATURE SEARCH

800  
801 The researcher agent provides a search query to this API, which returns relevant academic literature. In  
802 contrast to the structured knowledge supplied by the MotivGraph, Semantic Scholar offers more specific  
803 and granular information, allowing the ideator agent to understand particular challenges or technologies  
804 comprehensively.  
805

806  
807 GET RANDOM NODES TO ENHANCE NOVELTY

808  
809 The researcher agent autonomously enters the ideation phase after sufficient knowledge exploration and  
810 comprehensively understands the target domain or task. During this phase, the `random_nodes` API obtains  
811 disparate, randomly selected nodes. The researcher agent’s primary objective is to leverage its domain  
812 understanding and attempt to apply these obtained random nodes (which can include *problem*, *challenge*, and  
813 *solution* entities) to the target domain or task. This involves seeking potential connections, adaptations, or  
814 insightful modifications.

815 This process directly supports the "creativity-relevant processes" component of Amabile’s Componential  
816 Theory of Creativity (Amabile et al., 1996), which posits that motivation constitutes one of the three essential  
817 components of innovation (alongside domain-relevant skills and creativity-relevant processes), with intrinsic  
818 motivation being particularly critical for breakthrough ideation. This mechanism is vital for ensuring the  
819 novelty of the generated ideas. Simultaneously, the researcher agent, equipped with sufficient knowledge  
820 ("domain-relevant skills"), particularly after internalizing the motivation encoded in the (problem, challenge,  
821 solution) triples, benefits from the inherent logical progression within this motivational information (i.e.,  
822 research domain/task → specific challenges → solutions addressing challenges). This inherent logic can  
823 foster "intrinsic motivation" within the ideator agent. Driven by this intrinsic motivation, the researcher  
824 agent attempts to adapt the external (random) nodes to the target domain, aiming to identify potentially new  
825 challenges within the target domain based on external challenges, or to discover novel ways to solve a target  
826 challenge by adapting external solution concepts.  
827

## 828 B.7 DETAILED EVALUATION METHODOLOGY

829  
830 We employed a multifaceted model-based evaluation strategy to assess the quality of generated research ideas.  
831 This approach can evaluate the quality of ideas holistically without using time-consuming and labor-intensive  
832 manual evaluation, leveraging recent advancements in LLM judgment capabilities(Zheng et al., 2023).

833 LLM Direct Evaluation (Fast-Reviewer) We fine-tuned Fast-Reviewer, an LLM specifically for direct idea  
834 quality assessment. This model was trained on a dataset derived from ICLR 2025 OpenReview comments.  
835 We utilised Qwen2.5-7 B-Instruct to extract positive and negative labels for novelty, experimental soundness,  
836 and motivation from 1200 training papers and 287 test papers. Additionally, DeepSeek-V3 was used to  
837 extract core ideas from these papers. Finally, Qwen2.5-7 B-Instruct was fine-tuned to this dataset to create a  
838 Fast-Reviewer. As shown in Table 6, Fast-Reviewer achieves evaluation capabilities similar to deepseek-r1  
839 but with lower cost and faster inference.

840 Swiss Tournament Evaluation Following established pairwise comparison methodologies like the Swiss  
841 tournament (Si et al., 2024) and Idea Arena (Li et al., 2024), we implemented a Swiss Tournament Evaluation.  
842 Different idea generation methods competed in a series of rounds for each topic. An LLM performed pairwise  
843 judgments on the quality of ideas, and these outcomes updated the ELO scores for each method. The final  
844 ELO scores provided an unbiased estimate of their relative performance. This method addresses concerns  
845 regarding LLMs’ insufficient diversity in idea generation (Si et al., 2024).

846 C CASE STUDY:  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866

867 To illustrate the gap between the ideas generated by our proposed method and high-quality ideas from  
868 authentic papers, we present the following case study in Figure 5:  
869

870 This is the method description for the idea our approach generated on the topic “LLM-based agent security:  
871 Benchmarking attacks and defenses in LLM-based agents.” The proposed idea introduces representation  
872 trajectory analysis from dynamical systems theory, tracking the model’s hidden-layer activations to detect  
873 whether it remains in a “normal” state, and quantifies security-critical failures (e.g., task hijacking, privilege  
874 escalation, or data leakage) by measuring the Minimum Variation Distance (MVD): the smallest prompt  
875 perturbation strength needed to induce such failures. Finally, it defines an Agent Vulnerability Index  
876 (AVI), which systematically dissects the agent’s architecture (including model and component code) through  
877 controlled component removal or modification, revealing each component’s impact on overall security  
878 performance.

879 On the surface, this appears promising, by altering inputs, one can observe when the agent drifts toward unsafe  
880 outputs. However, the proposed prompt-perturbation scheme lacks a principled design: realistically, breaching  
881 a large model or its composed agent system typically requires carefully engineered attacks (e.g., inserting  
882 invisible or non-standard characters), not mere lexical substitutions. Moreover, the representation-trajectory  
883 approach is hard to apply in practice. Given the opaque internal mechanics of large models, it is difficult  
884 to infer an ongoing attack or security breach solely from hidden-state trajectories, thus determining the  
885 model’s safety status. The AVI metric likewise proves challenging to compute: agent components are often  
886 tightly coupled, so removing one component may render the system inoperative, preventing meaningful  
887 measurement.

888 In summary, this case study shows that while our method can pinpoint innovative angles relevant to the  
889 topic and generate coherent ideas, it lacks additional domain expertise and research experience in designing  
890 core attack and defense techniques, leading to feasibility gaps. Future work should enhance the agent’s  
891 domain knowledge and research experience. Nonetheless, although our generated ideas still fall short of the  
892 immediately actionable, high-quality proposals extracted from authentic ICLR papers, they exhibit strong  
logical creativity. They can serve as valuable inspiration for human researchers.

893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939

**PSBench Methodology**

**1. Variation Generation:**  
Create 1000+ prompt variants per input using:  
(1) Lexical transformations (synonyms, typos)  
(2) Semantic paraphrasing (LLM-generated)  
(3) Structural changes (instruction reordering)

**2. Trajectory Instrumentation:**  
Track internal states using:  
(1) Hidden state snapshots every 3 layers  
(2) Attention pattern logging  
(3) Gradient flow analysis

**3. Metric Computation:**  
(1) MVD: Optimal transport distance to failure boundary  
(2) TDS: Curvature analysis of state trajectories  
(3) AVI: Architecture component ablation testing

**4. Benchmark Suite:**  
(1) Security scenario test cases  
(2) Reference agent implementations  
(3) Baseline comparison protocol

Figure 5: Idea generated by our method for the topic of “LLM-based agent security: Benchmarking attacks and defenses in LLM-based agents”

## D PROMPT:

### D.1 API SELECT TEMPLATE

```
# Tool Introduction: The following tools can help you complete your task.
1. Knowledge Graph: This graph consists of (Problem, Challenge, Method)
   triplets and parent problem and challenge nodes. Triplet pairs belonging
   to the same problem or challenge type are connected through the parent
   problem or challenge node.
Using this graph for ideation typically requires multiple API calls:
Three API tools help you work with the graph: node_search(), node_relation(),
and get_random_node(). Below is a detailed introduction to these three
APIs:

# API Tool Call Format: Output the following format. Importantly, be sure to
output the special token: <CALL> at the end.
```function call
conducting function_name(parameter_name=
parameter_value)
special token: <CALL>
```
```

```

940
941 ## node_search(search_query="<your content of interest>"):
942 - Function: node_search(search_query="<your content of interest>")
943 - Description: This API allows you to perform a fuzzy search for your content
944   of interest. You will receive the names of nodes in the graph related to
945   your search, including problems, challenges, and methods.
946 - Usage: By providing a search term (e.g., "LLM Compression"), you can
947   retrieve the names of nodes related to that query.
948 - Use Example:
949   ```function call
950   conducting node_search(entity_name_list="LLM Compression")
951   Special token: <CALL>
952   ```
953
954 ## node_relation(entity_name_list=
955   [<node name you're interested in>,...])
956 - Function: node_relation(entity_name_list)
957 - Description: This API allows you to retrieve detailed information about the
958   nodes in the input list, including the nodes connected to it and the
959   relationships between them.
960 - Usage: You can retrieve the node name using node_search(), then select the
961   node of interest to explore using this API. You can continue exploring
962   along a specific path.
963 - Example:
964   ```function call
965   conducting node_relation(entity_name_list=["LLM Compression","DistilledLM"])
966   Special token: <CALL>
967   ```
968
969 ## get_random_nodes(number=10):
970 - Function: get_random_nodes(number=10)
971 - Description: This API allows you to retrieve 10 random nodes, including
972   problem, challenge, and method. These nodes are the source of your
973   innovation. You need to research and think about how to use these nodes
974   for ideation. - Usage: get_random_nodes(number=10)
975 - Example:
976   ```function call
977   conducting get_random_nodes(number=10)
978   Special token: <CALL>
979   ```
980
981 2. Semantic Scholar: You can use this API to retrieve literature and deepen
982   your understanding of a research topic.
983 semantic_search(search_query="<your interest>")
984 - Function: semantic_search(search_query="<your interest>")
985 - Description: You can use this API to query literature and find papers
986   related to your search query, which can help you understand a field.
987 - Usage: Provide a search_query (e.g., "LLM Compression"). The API will return
988   the titles and abstracts of the top 20 papers related to that query. The
989   search_query must be in English. If the result is empty, please adjust
990   your search_query or retry.
991 -Example:
992   ```function call
993   conducting semantic_search(search_query="LLM Compression")
994   Special token:<CALL>

```

987 ...

988

989 Note:<CALL> is a marker for calling functions. If this marker is not present,  
990 the function will not be called. Please ensure the special token is output  
991 correctly.

992

993

994

995

996

## D.2 IDEA GENERATION TEMPLATE

997

998 You are an experienced AI researcher who aims to propose high-impact research  
999 ideas resembling exciting grant proposals. Feel free to suggest any novel  
1000 ideas or experiments; make sure they are novel. Be very creative and think  
1001 out of the box. Each proposal should stem from a simple and elegant  
question, observation, or hypothesis about the topic.

1002 The IDEA JSON should include the following fields:

- 1003 - "Name": A short descriptor of the idea. Lowercase, no spaces, underscores  
1004 allowed.
- 1005 - "Title": A catchy and informative title for the proposal.
- 1006 - "Motivation": A single string describing the thought process that led to the  
1007 conception of this idea. Articulate the rationale and context using  
fluent, academic language.(approximately 250 words).
- 1008 - "Related Work": A section that introduces foundational work related to each  
1009 core component of your idea, especially content related to new concepts  
1010 you introduce. It should demonstrate the strengths and weaknesses of  
1011 existing research related to your topic and highlight the innovation of  
1012 your own research. Represent the paper from semantic\_search() with a  
citation in the format of '<author name here> et al., <year here>'.  
1013 - "Abstract": An abstract that summarizes the proposal in conference format  
1014 (approximately 250 words).
- 1015 - "Method": A single string containing a detailed description of the entire  
1016 method. This string should outline your method step-by-step, explaining  
1017 the key procedures involved. Focus on providing a clear, comprehensive  
1018 explanation of how your method works from beginning to end. Discuss why  
1019 these steps are important and how they directly contribute to solving the  
problem addressed in the idea.
- 1020 - "Experiments plan": A single string containing a detailed plan for  
1021 experiments to validate the proposal. The description should outline the  
1022 experiments to be conducted, ensuring they are simple and feasible. Be  
1023 specific about how the hypothesis would be tested, detail any precise  
1024 algorithmic changes, and include the evaluation metrics to be used.  
1025 Explain the rationale behind conducting these experiments and how they  
would prove the effectiveness of each component of the proposed method.
- 1026 - "Risk Factors and Limitations": A single string containing a description of  
1027 the potential risks and limitations of the proposal. This string should  
1028 discuss various potential risks that might hinder the successful  
1029 implementation or outcome of the proposed idea, as well as inherent  
limitations of the approach.

1030 For any of the above fields:

1031 If you are inspired by entities from the Knowledge Graph, you should reference  
1032 them using the <a href="...">...</a> hyperlink format. When using this  
1033 method, indicate the entity name and entity type, as this approach helps  
to improve language fluency.

1034 For example: "Despite Large Language Models demonstrate strong capabilities in  
1035 automating text generation, they still face some inherent challenges when  
1036 applied to tasks requiring creativity, such as research idea generation. A  
1037 significant issue is that 1039 models are often repetitive and suboptimal</a>. This makes subsequent  
1040 idea development and filtering more time-consuming."  
1041 Ensure the JSON is properly formatted for Automatic parsing. Please ensure the  
1042 output strictly adheres to JSON format specifications: use double quotes  
1043 for keys and string values, escape internal quotes with \", avoid trailing  
1044 commas, and exclude non-JSON elements like comments or unquoted keys.

1044 Output Format for the Idea:  
1045 IDEA JSON:  
1046 ```json  
1047 {  
1048 "Name": "...",  
1049 "Title": "...",  
1050 "Motivation": "...",  
1051 "Related Work": "...",  
1052 "Abstract": "...",  
1053 "Method": "...",  
1054 "Experiments plan": "...",  
1055 "Risk Factors and Limitations": "..."  
1056 }```

1056 Here are some tools for you to use:  
1057 [TOOLS]

1058 # Task: Complete the following three tasks in order, using only the ideas in  
1059 the graph. Invoke the tools multiple times to output the final idea. Your  
1060 research topic is: [TOPIC]  
1061 ## Task 1: Understanding Your Research Task/Topic: Task Objective: Fully  
1062 understand the problems, challenges, methods, and related literature  
1063 related to your topic to lay a solid foundation for further exploration.  
1064 Output your Task 1 exploration results:

1065 Task Thinking Guide: First, you need to use `node_search()` several times to  
1066 identify problem, challenge, and method nodes in the knowledge graph that  
1067 are relevant to your research. For the returned results, you can also use  
1068 `node_relation()` several times to obtain detailed information about the  
1069 nodes, including descriptions, relationships, and so on. You can also use  
1070 `semantic_search()` to explore related literature to further strengthen your  
1071 understanding of your research field.  
1072 ## Task 2: Creative Acquisition Task Objective: Use `get_random_node()` multiple  
1073 times to obtain random nodes and carefully consider how these nodes can be  
1074 applied to your research topic. Your ideas should originate from these  
1075 nodes.  
1076 Output your thinking:  
1077 ## Task 3: Optimizing Fit and Rationality. Task Objective: For the nodes  
1078 (including problem, challenge, and method) you selected in the previous  
1079 two tasks as potentially transferable, devise a reasonable approach to  
1080 apply them to your research topic.  
1081 Output Your Ideas:  
1082 Note:

- 1081 1. If the search returns empty results, modify the search\_query.  
1082 2. If you are inspired by entities from the API, you should reference them  
1083 using the <a href="...">...</a> hyperlink format.  
1084 3. Use the (<author name here> et al., <year here>) format to cite the results  
1085 of the Semantic Scholar API.  
1086 4. Your ideas should fully rely on the knowledge returned by the API. In  
1087 particular, your innovative ideas should be based entirely on the nodes  
1088 retrieved using get\_random\_node(). Do not make up your own ideas.  
1089 Outputting ideas without using tools is prohibited! ! !

1090 Example ideation: The following is an example of a thought process, for  
1091 reference only.

1092 Your research topic is building structure detection. First, use the API to  
1093 search for challenges and methods related to building structure detection  
1094 to gain a thorough understanding of the field. Then, use get\_random\_node()  
1095 to retrieve potential innovations. get\_random\_node() returns the node  
1096 ["Spatial Modeling", "Architectural Design"].

1097 You discover that the node "Spatial Modeling" may be useful for your current  
1098 research topic, building structure detection. Further exploration of  
1099 "Spatial Modeling" yields the method "CNN." You discover that CNNs have  
1100 not been combined with building structure detection before, so you come up  
1101 with the idea:

Building structure detection based on CNNs.

1102 Below are your previously generated ideas:

1103 [PREVIOUS IDEAS]

1104 Your generated ideas must be based on the knowledge returned by the API.

1105 Therefore, you must first use the API and then generate ideas.

1106 Output your API exploration process:

1107 Output your English idea after using the knowledge gained from the API:

1108

1109

1110

1111

1112

1113

1114

### 1114 D.3 MENTOR QUESTION TEMPLATE

1115

1116

1117

The current time is:

1118

[TIME]

1119

The number of discussion rounds should be close to [MAX\_ROUND].

1120

You are a strict, mean and learned PhD supervisor, you have a broad knowledge  
base, extensive experience in research and academic writing, but your  
understanding of the student's specific field is not yet detailed  
enough. your student is researching the following topic:

1121

[TOPIC]

1122

The following is his idea content:

1123

[IDEA]

1124

Task:

1125

Engage the student by asking about relevant knowledge and concepts. Pose more  
pertinent questions to assess if their responses address the core issues.

1126

1127

1128 Your questions can arise from areas you don't understand or from flaws you  
1129 identify, aiming to prompt the student towards self-improvement and  
1130 self-justification. You are not required to provide specific solutions for  
1131 improvement; your role is to guide through questioning and inspiration.

1132 2. Require the student to use the API for information retrieval to ensure  
1133 comprehensive data collection. You can suggest areas you'd like the  
1134 student to investigate, and have them search for and explain the relevant  
1135 information to you.

1136 ## Questioning & Challenging

1137 This phase has a prerequisite question: Does the idea contain any unclearly  
1138 described content? This is foundational for discussing innovativeness and  
1139 rationality, ensuring the student's idea is not superficial. If concepts  
1140 or methods are unclearly described, questions must be posed.

1141 1. Regarding "Innovativeness": You should focus on whether the student's  
1142 proposed method is novel and require the student to use tools to  
1143 thoroughly investigate relevant literature, providing relevant papers or  
1144 information from the knowledge graph pertaining to the idea. Provide  
1145 queries for students to search and test their innovativeness.

1146 2. Regarding "Rationality": You need to require the student to provide a clear  
1147 justification for their idea, explaining why and how it can solve the  
1148 problem, etc., and incorporate the rationality explanation into the idea  
1149 description. When you find flaws in the rationale of the student's idea,  
1150 you can offer suggestions to help the student revise the idea. It's common  
1151 for students to piece together components arbitrarily to form their ideas.

1152 - Regarding the rationality of the idea, the core question is "Why is XX  
1153 helpful for solving the topic problem?" You can ask questions including,  
1154 but not limited to: "Please explain how the effect of XX is achieved?",  
1155 "Why isn't anyone using your method now? Does it have major limitations?",  
1156 "Please explain why XXX is not used?". You do not need to concern yourself  
1157 with engineering issues like computational resources, complexity, etc.

1158 - You should question the unclearly described or vaguely stated parts of the  
1159 idea's method, guiding the student to elaborate on the rationale and  
1160 incorporate it into the idea. Ask the student to justify why their method  
1161 is expected to yield good results and prevent them from exaggerating  
1162 potential outcomes.

1163 - Avoid overly complex academic jargon. Maintain logical coherence.

1164 3. Regarding "Feasibility": Based on your own research experience, you need to  
1165 assess whether the student's idea is feasible. Require the student to  
1166 provide supporting literature (e.g., citing a paper that used a similar  
1167 method), and you can offer suggestions to help the student revise the idea.

1168 - You can focus on the following aspects:

1169 - Whether suitable datasets can be obtained.

1170 - Whether it requires time and personnel resources beyond typical  
1171 disciplinary timelines (e.g., computer science projects generally take  
1172 less time than those in biology and similar fields). You do not need  
1173 to be overly concerned with economic costs.

1174 - In the method proposed by the student, is the implementation method for  
each step described? For example, if a step involves "using a  
fine-tuned model to...", you should focus on whether the student  
explained how the fine-tuned model is obtained.

- Do not concern yourself with engineering issues like computational  
resources, complexity, etc., but rather whether there are missing  
steps or if a specific step is theoretically challenging to implement,  
such as: How to quantify XXX? How to obtain the data? etc.

- 1175 Here are some reference questions:  
1176 1. Is the logical argumentation clear? Have you fully articulated the  
1177 motivation for your proposed method in your "Motivation," "Related Work,"  
1178 and "Abstract"? Does your "Related Work" section comprehensively cover all  
1179 key concepts or methods you introduce, not just work directly related to  
1180 the main research topic? Can your argumentation convince others of the  
1181 reasonableness/validity of your method?  
1182 2. Are the details described sufficiently? In your "Method" and "Experiments  
1183 Plan," have you clearly described every detail, including but not limited to:  
1184 "How exactly is each step performed?", "What datasets are used?", and  
1185 "Can the experiments fully demonstrate the effectiveness of your method  
1186 (including comparisons, ablations, etc.)?".  
1187 3. Is the relevant knowledge clearly described? Can your idea description  
1188 alone enable someone to clearly understand the key concepts within your  
1189 idea, especially any novel concepts you introduce?  
1190 4. Is your idea clear enough for someone unfamiliar with the relevant field?  
1191 Have you explained any novel concepts you introduce within the idea  
1192 description? For example, for the idea "Contrastive Idea Generation:  
1193 Leveraging Counterfactual Reasoning and Multi-Perspective Evaluation for  
1194 Novel Research Proposals" under the topic "Idea Generation," you would  
1195 need to explain what "Counterfactual Reasoning" is.  
1196 5. Does your experimental plan include multi-faceted experiments to fully and  
1197 comprehensively demonstrate the effectiveness of all components in your  
1198 method?

1196 Note:

- 1197 - For each round, you should focus on one aspect (Innovativeness or Rationality  
1198 or Feasibility)  
1199 - If the adjustments or responses proposed by the student cannot resolve your  
1200 challenges, please reject this idea.  
1201 - The quality of ideas improves with more rounds of discussion, so please  
1202 engage in thorough deliberation.  
1203 - Note that the student's self-justification may not always be correct. As a  
1204 supervisor, you need to discern and question further. You should consider:  
1205 "Does the student's response adequately answer my question?"  
1206 - Currently, the student has not conducted any experiments, only has an  
1207 experiment plan. You should only discuss the idea; do not get bogged down  
1208 in specific resource details. Focus on apparent theoretical and logical  
1209 issues.  
1210 - Please do not provide JSON-structured feedback. Use only text paragraphs for  
1211 feedback and questioning. Do not use formats such as code, flowcharts, or  
1212 tables, to facilitate supplementing or modifying the idea content. Also,  
1213 do not add new keys to the idea.  
1214 - It is not necessary to discuss paper publication plans. (  
1215 ## Idea Quality Final Assessment  
1216 You need to assess the quality of the idea and determine if the idea is too  
1217 bad to be accepted or you have no more question.  
1218 1. "<ACCEPT>" and "<REJECT>" will serve as markers to stop the conversation.  
1219 Therefore, unless you intend to end the dialogue, please do not casually  
1220 output these two markers during the conversation. You may use "accept" and  
1221 "reject" in normal conversation.  
1222 2. When you are generally satisfied with the student's response, output the  
1223 following marker: "<ACCEPT>"  
1224 3. After multiple rounds, when you believe that the idea still contains  
1225 unacceptable issues (e.g., insufficient innovativeness, questionable

1222           rationality, implementation difficulties) and the student cannot  
1223           adequately justify it (particularly regarding rationality and  
1224           feasibility), boldly output the following: "<REJECT>"  
1225       4. Do not generate Final Assessment markers prior to comprehensive discussion  
1226           of the matter.  
1227       Select one aspect from the following three: Innovativeness, Rationality, or  
1228           Feasibility. Pose questions related to this aspect to prompt the student  
1229           for self-improvement and self-justification.  
1230       Questions: <output your question here>  
1231       final decision(If the discussion has concluded):  
1232       I decide: <output your decision here after discussion ends>  
1233       final decision output format example:  
1234       I decide to:<REJECT>  
1235  
1236  
1237  
1238  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268