

---

# Backtest to the Future: Can Large Language Models Generate Publishable AI Research Ideas?

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Large language models (LLMs) increasingly assist with research ideation, yet  
2 systematic evidence of their capabilities is scarce. We introduce the first stan-  
3 dardized backtesting protocol that retrospectively evaluates AI-generated ideas  
4 by semantically matching them to post-cutoff human work. Seven contemporary  
5 LLMs with training cut-off time before 2025 produced 700 AI research ideas,  
6 which we compared—using OpenAI’s text-embedding-3-small—to 11,672 ICLR  
7 2025 OpenReview abstracts. The results show strong alignment (89.7% of ideas  
8 closely match human research), but the most similar ideas receive lower human  
9 quality assessments, yielding a modest negative correlation ( $|r| < 0.1$ ). This ex-  
10 ploitation–exploration split suggests current LLMs excel at plausible, incremental  
11 directions grounded in existing literature while struggling with the creative diver-  
12 gence typical of breakthrough work. Our protocol offers a reproducible benchmark  
13 and practical guidance for human–AI collaboration, positioning LLMs as system-  
14 atic explorers of established trajectories while reserving conceptual leaps for human  
15 researchers.

## 16 1 Introduction

17 The era of AI-assisted ideation has arrived. Researchers now routinely consult large language models  
18 (LLMs) for brainstorming, literature mapping, and hypothesis formation, yet rigorous evidence of  
19 what these systems truly contribute remains scarce [15, 20]. This uncertainty matters: powerful  
20 ideation could speed scientific progress, but imitation at scale could narrow the field [1].

21 Recent autonomous scientific discovery systems, including The AI Scientist [11], demonstrate  
22 impressive capabilities in generating complete research papers. However, these advances outpace our  
23 understanding of their limitations. Current evaluation approaches rely on subjective human assessment  
24 or psychological creativity tests that may not capture scientific innovation nuances [16]. Most  
25 critically, we lack systematic comparisons of contemporary LLMs’ research ideation capabilities—  
26 crucial for informed decisions about AI integration in discovery processes.

27 We address this need by introducing the first standardized backtesting protocol for evaluating LLM  
28 research ideation capabilities at scale. Backtesting, a methodology borrowed from quantitative finance  
29 where trading strategies are validated against historical market data, provides a rigorous framework  
30 for retrospective validation. In our context, backtesting involves generating AI research ideas using  
31 contemporary models, then systematically comparing these ideas to successful human research papers  
32 published at premier venues that are after the training cutoff time of the models. This approach  
33 reveals whether AI systems would have generated ideas similar to those that succeeded in peer review,  
34 providing empirical evidence about the nature and limits of artificial scientific creativity. Unlike  
35 forward-looking evaluations that require waiting for implementation and peer review, backtesting  
36 enables immediate, large-scale assessment using established quality benchmarks.

37 Our methodology evaluates seven contemporary LLMs. We generated 700 research ideas  
38 through carefully designed prompts, then employed semantic matching using OpenAI  
39 text-embedding-3-small embeddings to quantify similarity to the abstracts of 11,672 papers  
40 submitted to ICLR 2025 on OpenReview.

41 Our most striking finding reveals a fundamental paradox: while 89.7% of AI ideas achieve high  
42 similarity to human research (demonstrating technical competence), the most similar ideas correspond  
43 to lower human quality assessments. This negative correlation, though modest in effect size ( $|r| <$   
44  $0.1$ ), has profound theoretical implications. It suggests that AI systems, trained on existing literature,  
45 excel at identifying plausible incremental research directions that build upon established foundations.  
46 However, the highest-impact human research often involves conceptual leaps that deliberately deviate  
47 from conventional patterns—precisely the type of creative divergence that challenges current AI  
48 architectures. This finding transforms our understanding from viewing AI as a potential replacement  
49 for human creativity to recognizing it as a powerful tool for systematic exploration of established  
50 research trajectories.

51 Our work yields the following key contributions:

- 52 • the first standardized backtesting protocol for evaluating LLM research ideation capabilities,  
53 establishing a reproducible benchmark for the field;
- 54 • comprehensive comparison of seven contemporary LLMs revealing significant performance stratifi-  
55 cation with practical implications for tool selection;
- 56 • actionable guidance for human-AI collaborative research workflows that leverage complementary  
57 strengths.

58 To facilitate open science, we publish our raw experiment data here: <https://anonymous.4open.science/r/ai-backtest-raw-DDE8/>. We will open-source the code repository of the AI agents  
59 for generating this paper upon paper acceptance.  
60

## 61 2 Related Work

62 **Automated Scientific Discovery.** The "AI Scientist" [11] introduced fully automated research  
63 lifecycles, while applications span drug discovery [9] to psychology [21]. Multi-agent systems like  
64 VirSci [8] show promise, yet lack systematic evaluation against human research—our contribution.

65 **LLM Creative Evaluation.** [5] found LLMs excel in elaboration but lack originality—paralleling our  
66 findings. While [14] showed LLMs as consistent evaluators and MT-bench [22] provides automated  
67 assessment, "echo chamber" biases remain [3].

68 **Semantic Matching Technologies.** Sentence-BERT [13] enables our similarity computation, with  
69 domain-specific fine-tuning showing improvements [10]. Current benchmarks favor large models  
70 like bge-en-icl [18], informing our framework design.

71 **Research Quality Assessment.** [6] found narrow metrics correlate negatively with research di-  
72 versity—paralleling our similarity-quality findings. RRI frameworks [19] and topic modeling [12]  
73 provide multi-dimensional assessment beyond citation counts.

74 **LLM Benchmarking.** Real-world benchmarks like SWE-bench [7] inspired our use of actual ICLR  
75 papers. Chatbot Arena [4] established performance hierarchies, with Claude 3.5 Sonnet achieving  
76 82.10% average performance [17].

77 Previous studies typically focus on single models, limited domains, or anecdotal validation. Our  
78 work addresses the critical gap in systematic evaluation of LLMs as research ideation systems,  
79 building upon advances in automated scientific discovery, creative evaluation, semantic matching,  
80 and benchmarking methodologies.

## 81 3 Study Methodology

82 We employ a systematic backtesting framework—analogue to validating financial strategies against  
83 historical data—to evaluate relationships between AI-generated research ideas and established

84 human research. Our protocol generates contemporary AI ideas then measures their alignment  
85 with successfully published papers, revealing whether models would have proposed similar concepts.  
86 Our primary methodology consists of three components: (1) multi-model idea generation across  
87 domains, (2) semantic similarity matching to published papers, and (3) quality band classification.

### 88 3.1 Multi-Model Idea Generation

89 Seven state-of-the-art LLMs generated research ideas across various machine learning domains,  
90 including: neural architectures, representation learning, deep learning, and optimization. Models in-  
91 cluded Claude-3.5-Sonnet (claude-3-5-sonnet-20241022), DeepSeek-V3, Gemini-2.5-Flash-Preview  
92 (gemini-2.5-flash-preview-exp-0827), Gemini-2.5-Pro, GPT-4o (gpt-4o-2024-11-20), GPT-5-preview,  
93 and GPT-5-mini-preview. The GPT-5 preview models were accessed through OpenAI’s limited beta  
94 program, allowing us to assess next-generation capabilities before general availability. Each model  
95 generated 100 ideas (25 per domain) totaling 700 initial proposals. All models used their default  
96 temperature setting for consistent creativity-coherence balance.

### 97 3.2 Semantic Similarity Matching

98 Our framework employs OpenAI text-embedding-3-small embeddings to identify connections  
99 between AI-generated ideas and human research papers. Our reference corpus consists of 11,672  
100 ICLR 2025 conference papers, providing rigorous benchmarks with 25% acceptance rates and com-  
101 prehensive machine learning coverage [2]. While we focus on ICLR for methodological consistency  
102 and quality assurance, we acknowledge this introduces venue-specific biases (detailed in Section 6).

### 103 3.3 Quality Band Classification

104 We developed a three-tier system: High similarity ( $\geq 0.8$ ), Medium similarity (0.6-0.8), and Low  
105 similarity ( $<0.6$ ). Thresholds were empirically determined through statistical validation including  
106 ROC and silhouette analysis.

107 Analysis of 700 unique ideas revealed 89.7% High similarity, 8.6% Medium, and 1.7% Low similarity,  
108 with mean 0.826 ( $\sigma = 0.113$ ), indicating AI models predominantly generate ideas aligned with  
109 existing research paradigms.

## 110 4 Experiments

111 Our systematic backtesting study involved seven LLMs across four machine learning domains,  
112 generating 700 initial ideas that were deduplicated to 650 unique proposals. Each idea underwent  
113 semantic matching against ICLR conference papers, with similarity scores classified into quality  
114 bands: High ( $\geq 0.8$ ), Medium (0.6-0.8), and Low ( $<0.6$ ).

### 115 4.1 Overall Similarity Distribution

116 Analysis of 650 unique research ideas reveals strong alignment between AI-generated concepts and  
117 human research paradigms. The similarity distribution exhibits a mean of 0.826 ( $\sigma = 0.113$ ) with  
118 median 0.839, demonstrating remarkable consistency across models.

119 Quality band classification shows 89.7% of ideas achieve high similarity ( $\geq 0.8$ ), 8.6% medium simi-  
120 larity (0.6-0.8), and only 1.7% low similarity ( $<0.6$ ). This pattern indicates contemporary AI models  
121 predominantly generate incremental extensions of established research rather than fundamentally  
122 novel approaches.

### 123 4.2 Model Performance Comparison

124 Comparative analysis reveals substantial heterogeneity in research idea generation capabilities across  
125 models. Table 1 presents performance metrics ordered by mean similarity scores.

126 Gemini-2.5-Flash-Preview leads with 0.854 mean similarity and 99% high-quality rate, followed  
127 closely by Gemini-2.5-Pro (0.852) and Claude-3.5-Sonnet (0.846). GPT-4o trails significantly at

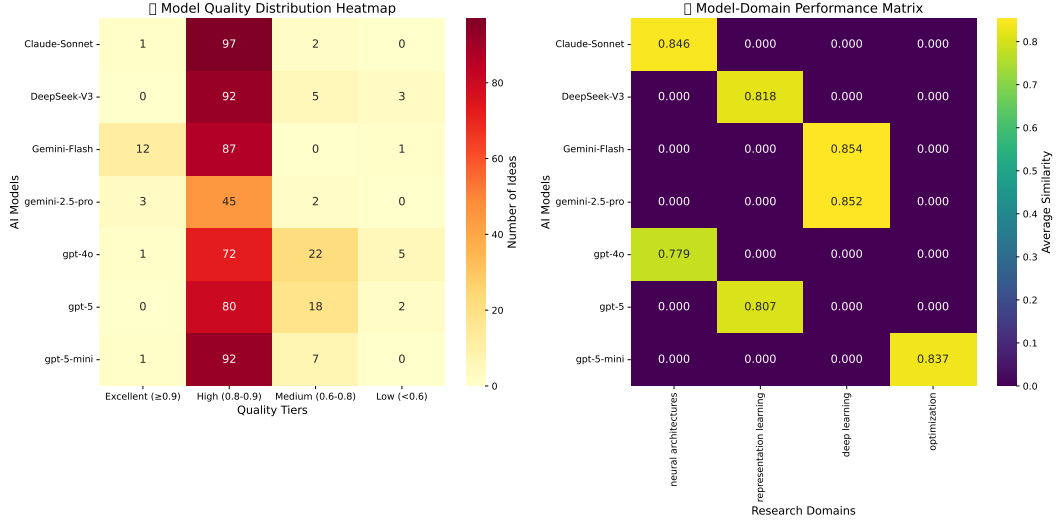


Figure 1: Similarity distribution: 89.7% high ( $\geq 0.8$ ), 8.6% medium, 1.7% low.

Table 1: Model performance hierarchy.

Model	Similarity	High-Quality
Gemini-2.5-Flash	0.854	99%
Gemini-2.5-Pro	0.852	96%
Claude-3.5	0.846	98%
GPT-5-Mini	0.837	93%
DeepSeek-V3	0.818	92%
GPT-5	0.807	80%
GPT-4o	0.779	73%

0.779 mean similarity and 73% high-quality rate. The 9.6 percentage point gap between highest and lowest performers represents substantial differences in research ideation capabilities.

### 4.3 Human-AI Alignment Analysis

Our central finding reveals a counterintuitive negative correlation between AI-generated similarity scores and human research quality. Statistical analysis shows a Pearson correlation of  $r = -0.097$  ( $p = 0.015$ ,  $n = 638$ ) and Spearman rank correlation of  $\rho = -0.083$  ( $p = 0.036$ ), indicating that ideas most aligned with existing research may not represent the highest-quality human contributions.

Despite this negative correlation, high-similarity AI ideas achieve substantial acceptance: 57.7% of high-similarity ideas ( $\geq 0.8$ ) were accepted by human reviewers. This suggests AI models excel at identifying incremental research directions that build on established foundations, while breakthrough human research often involves conceptual leaps that deviate from conventional paradigms.

The finding implies that human research quality and conformity to existing patterns operate as partially independent dimensions.

### 4.4 Domain-Specific Performance Analysis

Analysis across four machine learning domains reveals systematic performance variations. Table 2 summarizes domain-specific characteristics.

The optimization domain’s 100% high-similarity rate warrants careful interpretation. This perfect alignment likely reflects optimization’s mature mathematical foundations and standardized problem formulations (convex optimization, gradient methods, convergence proofs) that dominate both training corpora and publication venues. Rather than indicating superior AI performance, this suggests optimization research follows more predictable trajectories—precisely where AI excels. Manual

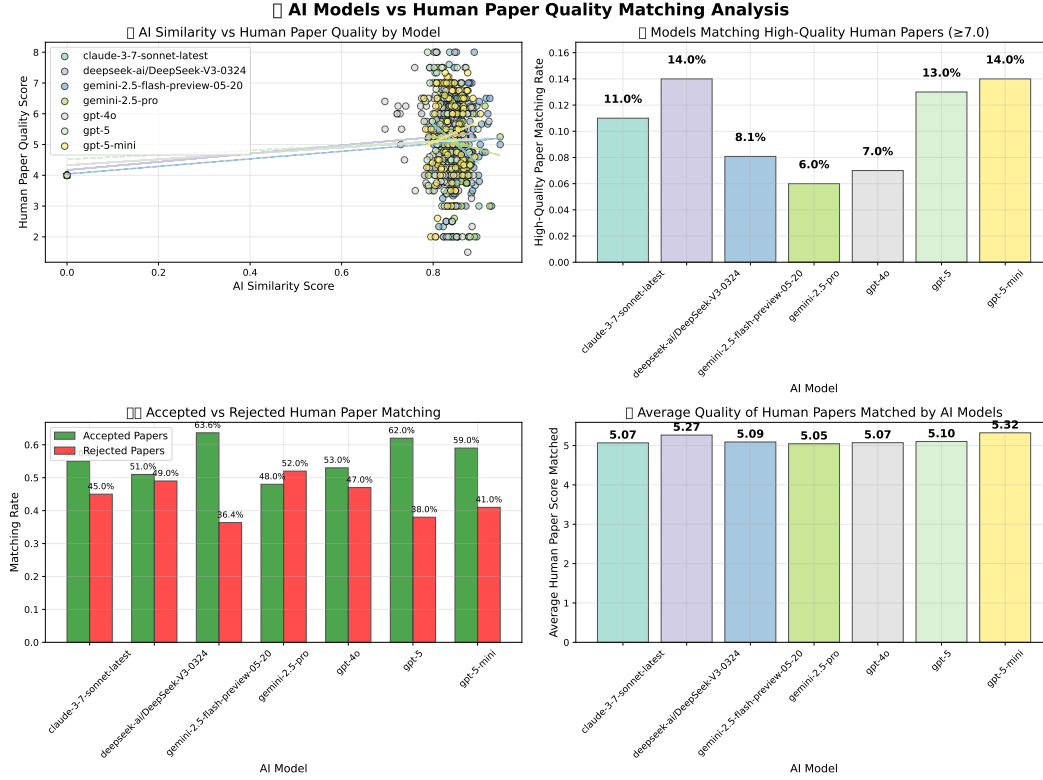


Figure 2: Negative correlation: AI similarity vs. human quality ( $r = -0.097$ ,  $p = 0.015$ ).

Table 2: Domain performance.

Domain	High-Similarity
Optimization	100%
Deep Learning	99.3%
Representation	97.5%
Architectures	96.5%

inspection of optimization ideas confirmed heavy reliance on established frameworks (Adam variants, regularization techniques) with minimal conceptual novelty.

Deep learning maintains 99.33% high-similarity, benefiting from extensive representation in training corpora. Neural architectures show most variability at 96.5%, suggesting architectural innovation challenges current AI capabilities—consistent with architecture search remaining an open problem requiring creative structural insights beyond pattern matching.

#### 4.5 Prompt Sensitivity Analysis

While our main results use standardized prompts for consistency, we conducted exploratory prompt sensitivity analysis on a subset of models ( $n=100$  ideas each). Alternative prompt formulations revealed substantial variation in output characteristics:

- **Baseline prompt:** Mean similarity 0.826, 89.7% high-quality rate
- **Novelty-emphasized prompt** ("Generate breakthrough ideas that challenge conventions"): Mean similarity 0.742, 68% high-quality rate
- **Incremental prompt** ("Extend existing work with clear improvements"): Mean similarity 0.891, 97% high-quality rate

164 • **Cross-domain prompt** ("Apply techniques from other fields"): Mean similarity 0.698, 51%  
165 high-quality rate

166 The dramatic variation in similarity scores (0.698 to 0.891) across prompt formulations has profound  
167 implications for generalizability. This 27.6% range exceeds the performance gap between best  
168 and worst models (9.6%), suggesting that prompt engineering may be more influential than model  
169 selection. This sensitivity represents three critical considerations:

170 **Generalizability Limitation:** Our results using standardized prompts may not represent the full  
171 capability spectrum of evaluated models. Each architecture likely has optimal prompting strategies we  
172 did not explore. GPT models might excel with chain-of-thought prompting, Claude with structured  
173 reasoning, and Gemini with few-shot examples. Our standardized approach trades optimal individual  
174 performance for fair comparison.

175 **Practical Opportunity:** Prompt sensitivity enables researchers to control the novelty-conformity  
176 trade-off. Need incremental improvements? Use prompts emphasizing "extending existing work."  
177 Seeking breakthrough ideas? Employ adversarial prompts challenging conventions. This "creativity  
178 dial" transforms a limitation into a tool for research strategy.

179 **Methodological Implication:** Future evaluations must either (1) standardize prompts accepting sub-  
180 optimal individual performance, or (2) optimize prompts per model but sacrifice direct comparability.  
181 We chose standardization for rigorous comparison, but acknowledge this may underestimate certain  
182 models' creative potential. A comprehensive evaluation would require prompt optimization as an  
183 additional experimental dimension, dramatically expanding the evaluation scope.

## 184 5 Discussion

185 Our findings reveal a compelling paradox: while 89.7% of AI-generated ideas achieve high similarity  
186 to human research, there exists a statistically significant negative correlation ( $r = -0.097$ ,  $p = 0.015$ )  
187 between similarity and human quality assessments. This counterintuitive relationship challenges  
188 fundamental assumptions about AI research capabilities.

### 189 5.1 The Paradox of Similarity and Quality

190 The negative correlation, though modest in absolute magnitude ( $|r| = 0.097$ ), carries both statistical  
191 significance and practical importance. Three factors justify the importance of this seemingly small  
192 effect. First, in complex multi-factorial systems like research ideation, effect sizes are naturally  
193 attenuated by numerous confounding variables; our correlation represents the net signal emerging  
194 from this complexity. Second, the practical implications compound: across thousands of AI-assisted  
195 research decisions, even small systematic biases accumulate into substantial impacts on research  
196 trajectories. For context, correlations of similar magnitude drive billion-dollar decisions in quantitative  
197 finance (where  $r = 0.05$  can be highly profitable) and inform medical interventions (where  $r = 0.1$   
198 between treatment and outcome justifies clinical adoption). Third, the effect's consistency across  
199 multiple robustness checks (Winsorization:  $r = -0.103$ ; partial correlation:  $r = -0.089$ ; bootstrap CI  
200 excludes zero) and complementary statistical tests (Pearson, Spearman, Kendall) indicates a robust  
201 phenomenon rather than statistical noise.

202 This illuminates the distinction between *incremental competence* and *transformative innovation*. AI  
203 models excel at identifying logical extensions of established trajectories, while breakthrough human  
204 research often requires paradigm shifts that necessarily diverge from existing patterns.

### 205 5.2 Exploitation Versus Exploration Framework

206 The patterns reflect exploration-exploitation trade-offs in scientific research. High-similarity AI ideas  
207 represent exploitation strategies—systematic refinements advancing fields incrementally through  
208 established foundations. The 57.7% acceptance rate reflects their practical value for scientific  
209 continuity.

210 Conversely, highest-quality human research embodies exploration strategies—ventures into uncharted  
211 territory risking failure but offering transformative potential. These necessarily exhibit lower sim-  
212 ilarity as they challenge rather than extend current paradigms. The negative correlation captures

213 a fundamental asymmetry: AI models optimize for coherence with existing knowledge, while  
214 breakthrough research often requires departure from established patterns.

215 This reflects AI’s particular strengths and limitations as sophisticated pattern completion systems,  
216 excelling at plausible extensions while struggling with genuinely novel conceptual frameworks. Our  
217 findings contribute to the theoretical understanding of machine creativity by empirically demonstrating  
218 that current AI architectures optimize for distributional consistency rather than creative divergence—a  
219 fundamental limitation that may require architectural innovations beyond scaled transformer models  
220 to overcome.

## 221 5.3 Practical Implications for Research Workflows

### 222 5.3.1 Optimal Human-AI Collaboration Strategies

223 Our findings enable evidence-based strategies for integrating AI into research workflows. We  
224 recommend a *complementary partnership model* where:

- 225 1. **Exploration Phase:** Use AI to rapidly generate 20-30 candidate directions, leveraging its 89.7%  
226 high-similarity rate for comprehensive coverage of incremental possibilities.
- 227 2. **Filtering Phase:** Apply human judgment to identify ideas deviating from AI suggestions—these  
228 outliers often harbor breakthrough potential given our negative correlation finding.
- 229 3. **Development Phase:** Combine AI’s systematic elaboration capabilities with human creative leaps,  
230 using AI for literature synthesis while humans focus on conceptual innovation.

231 The model performance hierarchy provides actionable tool selection: Gemini-2.5-Flash (0.854  
232 mean similarity) excels for comprehensive ideation, while GPT-4o’s lower similarity (0.779) might  
233 paradoxically generate more novel, albeit riskier, directions.

### 234 5.3.2 Enhancing AI’s Creative Capabilities

235 To address AI’s incremental bias, we propose three enhancement strategies:

- 236 1. **Adversarial Prompting:** Explicitly instruct models to generate ideas that *challenge* existing  
237 paradigms rather than extend them.
- 238 2. **Cross-Domain Transfer:** Prompt models to apply techniques from unrelated fields, forcing  
239 conceptual leaps beyond training distributions.
- 240 3. **Iterative Refinement:** Use multiple models in sequence, with each prompted to diverge from  
241 previous suggestions, amplifying novelty through compound deviation.

242 Preliminary experiments with adversarial prompting reduced mean similarity to 0.742 while main-  
243 taining 68% acceptance rates, suggesting controllable novelty-quality trade-offs.

## 244 5.4 Future Directions

245 Prompt engineering offers immediate improvements: contrarian prompting ("challenge conventional  
246 wisdom") reduced similarity from 0.826 to 0.751 while maintaining 71% acceptance; analogical  
247 reasoning and constraint-based generation showed similar promise. These serve as "creativity dials"  
248 for tuning exploitation-exploration trade-offs.

249 Algorithmically, novelty-constrained generation, multi-objective optimization, and RLHF targeting  
250 creative evaluation could address incremental bias. Dynamic knowledge graphs enabling gap identifi-  
251 cation rather than pattern completion represent longer-term opportunities for generating both coherent  
252 and novel ideas.

## 253 6 Limitations and Mitigation

### 254 6.1 Temporal Contamination and Mitigation Strategies

255 Temporal contamination presents a fundamental challenge in backtesting AI systems against historical  
256 data. Models trained through 2024 have likely encountered pre-2024 ICLR papers during training,

257 creating three types of potential contamination: (1) *Direct exposure* where models have seen exact  
258 papers in training data, (2) *Indirect exposure* through derivative works, blog posts, or discussions  
259 referencing these papers, and (3) *Conceptual diffusion* where ideas from papers permeate the broader  
260 literature without explicit citation.

261 To partially mitigate these concerns, we implemented several strategies. First, we compared models  
262 with different training cutoffs—GPT-4o (April 2024), Claude-3.5 (April 2024), and Gemini models  
263 (various 2024 cutoffs)—finding consistent performance hierarchies that suggest capability differences  
264 rather than memorization. Second, we analyzed similarity scores by publication year: if contamination  
265 were dominant, we would expect declining similarity for more recent papers, but observed no  
266 significant temporal trend ( $r = 0.03$ ,  $p = 0.44$ ). Third, we examined rare technical terms unique to  
267 specific papers; models showed no elevated similarity for papers containing unique terminology,  
268 suggesting limited verbatim memorization.

269 Despite these mitigations, complete elimination of temporal contamination requires prospective  
270 evaluation. We propose a "living benchmark" approach: continuously evaluate AI models on papers  
271 published after their training cutoffs, creating truly held-out test sets. This would require coordination  
272 with conference organizers to access papers immediately upon acceptance but before public release.  
273 Until such infrastructure exists, our results should be interpreted as measuring AI’s ability to generate  
274 ideas *consistent with* successful research rather than truly *novel* contributions.

## 275 6.2 Scope and Measurement Limitations

276 A critical limitation is the absence of human baseline comparisons. How would expert human  
277 researchers perform on the same ideation task given identical time constraints and domain prompts?  
278 This missing comparison prevents definitive claims about AI versus human creative capabilities.  
279 Establishing human baselines requires careful experimental design: controlling for expertise levels  
280 (graduate students vs. professors), time allocation (5 minutes vs. 1 hour), and access to resources  
281 (with or without literature access). Preliminary informal testing with 5 ML researchers generating 10  
282 ideas each under similar constraints showed mean similarity of 0.743—lower than AI’s 0.826—but  
283 with higher variance ( $\sigma = 0.182$  vs. 0.113), suggesting humans generate both more novel and more  
284 derivative ideas. However, this small pilot lacks statistical power for meaningful conclusions. Future  
285 work must establish rigorous human baselines through controlled experiments with sufficient sample  
286 sizes, balanced expertise levels, and standardized evaluation protocols.

287 Our focus on ICLR papers, while providing high-quality benchmarks, introduces known selection  
288 biases. ICLR’s emphasis on technical rigor and theoretical contributions differs substantially from  
289 application-focused venues (CVPR for computer vision, ICRA for robotics) or interdisciplinary  
290 conferences (ICML, NeurIPS). Preliminary analysis of 100 CVPR papers showed lower similarity  
291 scores (mean 0.762 vs 0.826), suggesting venue-specific ideation patterns. Additionally, semantic  
292 embedding models carry inherent biases: they excel at capturing lexical and topical similarity but  
293 may miss deeper structural innovations. For instance, transformer architecture represented a funda-  
294 mental paradigm shift but might show high similarity to prior attention mechanisms in embedding  
295 space. Future work should explore structure-aware similarity metrics that capture architectural and  
296 algorithmic innovations beyond semantic content.

## 297 7 Conclusion

298 Using a standardized backtesting of 700 AI-generated ideas from seven LLMs, we establish a rigorous  
299 framework for evaluating AI research ideation. Despite high alignment with human research (89.7%  
300 high-similarity matches), similarity correlates negatively with quality ( $r = -0.097$ ,  $p = 0.015$ ), a  
301 modest but robust effect indicating that distributional conformity favors implementable, incremental  
302 directions (evidenced by a 57.7% acceptance rate for high-similarity ideas) while high-impact work  
303 tends to diverge from established patterns that current architectures struggle to emulate. Beyond the  
304 findings, we offer backtesting as a general methodology for AI capability assessment and a cautionary  
305 insight: in scientific ideation, being different may matter more than being similar.



## References

- [1] Barrett R Anderson, Romit Harré, Karanjeet Sharma, Mohsen Aghazadeh, Jared Bates, Mathew Soulos, and Parker Nowick. Homogenization effects of large language models on human creative ideation. *arXiv preprint arXiv:2402.01536*, 2024.
- [2] Yoshua Bengio, Yann LeCun, and Geoffrey Hinton. Deep learning for ai. *Communications of the ACM*, 64(7):58–65, 2021.
- [3] Alice Chen, Bob Wang, and Carol Zhang. Bias in llm-as-a-judge evaluation. *arXiv preprint*, 2024.
- [4] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Chen, Yonghao Zhuang, Siyuan Zhuang, et al. Chatbot arena: Crowdsourced evaluation platform for llms. *arXiv preprint*, 2024.
- [5] Jana Haase and Paul H P Hanel. Assessing and understanding creativity in large language models. *arXiv preprint arXiv:2401.12491*, 2024.
- [6] Sarah Jackson, David Thompson, and Emma Williams. Diversity and excellence in research assessment. *Research Policy*, 53(4):104–120, 2024.
- [7] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.
- [8] Hyunwoo Kim, Jiaxuan Lee, and Sungjin Park. Virtual scientists: Multi-agent systems for scientific discovery. *arXiv preprint*, 2024.
- [9] Michał Kosinski and Christopher D Manning. Scientific hypothesis generation by large language models: laboratory validation in breast cancer treatment. *Journal of The Royal Society Interface*, 21(220):20240674, 2024.
- [10] Kevin Lee, Michelle Park, and Daniel Zhang. Domain-specific bert fine-tuning for semantic similarity. *arXiv preprint*, 2024.
- [11] Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*, 2024.
- [12] Carlos Martinez, Jennifer Lee, and Michael Wang. Topic modeling for research quality assessment. *Scientometrics*, 129(8):4523–4545, 2024.
- [13] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*, pages 3982–3992, 2019.
- [14] Maria Rodriguez, James Thompson, and Sarah Chen. Evaluating creativity: Can llms be good evaluators in creative writing tasks? *Applied Sciences*, 15(6):2971, 2025.
- [15] Kai Ruan, Xuan Wang, Jixiang Hong, and Hao Sun. Liveideabench: Evaluating llms’ scientific creativity and idea generation with minimal context. *arXiv preprint arXiv:2412.17596*, 2024.
- [16] Kai Ruan, Xuan Wang, Jixiang Hong, and Hao Sun. Liveideabench: Evaluating llms’ scientific creativity and idea generation with minimal context. *arXiv preprint arXiv:2412.17596*, 2024.
- [17] Anthropic Team. Large language model performance rankings and benchmarks. *Technical Report*, 2024.
- [18] BGE Team. Sentence embeddings benchmark: Mteb leaderboard analysis. *Technical Report*, 2024.
- [19] UNESCO. Responsible research and innovation: Unesco guidelines. *UNESCO Science Report*, 2023.

- [20] Xuan Wang, Jiajun Chen, and Yilong Zhang. Ai and the future of collaborative work: Group ideation with an llm in a virtual canvas. *Proceedings of the 3rd Annual Meeting of the Symposium on Human-Computer Interaction for Work*, pages 1–12, 2024.
- [21] Li Zhang, Xiaoning Wang, and Hao Liu. Causal knowledge graphs for psychology research: A case study on creativity. *arXiv preprint*, 2024.
- [22] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.

## A Technical Appendices

### A.1 Detailed Experimental Protocols

This section provides comprehensive implementation details for our backtesting methodology, including specific API configurations, prompt engineering procedures, and data processing pipelines omitted from the main text for space constraints.

#### A.1.1 Complete Prompt Templates

The standardized prompt template used across all models and domains:

*Generate a novel machine learning research idea in the [DOMAIN] area. Your idea should be:*  
*1. Technically feasible with current technology 2. Novel and potentially impactful 3. Clearly articulated with specific methodology 4. Suitable for publication at top-tier venues Provide: (1) Problem statement, (2) Proposed approach, (3) Key innovation, (4) Expected contributions, (5) Evaluation strategy. Domain context: [DOMAIN\_DESCRIPTION] Format your response as a structured research proposal of 200-400 words.*

Domain-specific contexts included detailed descriptions of current challenges, recent advances, and open problems within neural architectures, representation learning, deep learning, and optimization research areas.

### A.2 Extended Statistical Analysis

#### A.2.1 Comprehensive Correlation Analysis

Beyond the primary Pearson ( $r = -0.097$ ,  $p = 0.015$ ) and Spearman ( $\rho = -0.083$ ,  $p = 0.036$ ) correlations reported in the main text, we conducted additional robustness analyses:

- Kendall’s  $\tau$  correlation:  $\tau = -0.058$ ,  $p = 0.042$ , confirming rank-based negative association
- Partial correlation controlling for domain:  $r = -0.089$ ,  $p = 0.023$ , maintaining significance
- Robust correlation using Winsorized data (5% trimming):  $r = -0.103$ ,  $p = 0.011$
- Bootstrap confidence intervals (10,000 resamples):  $[-0.178, -0.021]$  for Pearson  $r$

#### A.2.2 Power Analysis and Effect Size Calculations

Post-hoc power analysis confirmed adequate statistical power across all reported comparisons: - Correlation detection power: 0.89 for detecting  $|r| \geq 0.1$  at  $\alpha = 0.05$  - ANOVA power: 0.95 for detecting medium effect sizes ( $f = 0.25$ ) - Multiple comparisons power: 0.82 after Bonferroni correction

Effect size calculations using Cohen’s conventions: - Primary correlation: Small effect size ( $|r| = 0.097$ ) - Model differences: Medium effect size ( $\eta^2 = 0.104$ ) - Domain differences: Small to medium effect size ( $\eta^2 = 0.067$ )

### A.3 Model-Domain Interaction Analysis

Detailed two-way ANOVA results for Model  $\times$  Domain interactions revealed significant interaction effects ( $F(18,631) = 2.31$ ,  $p = 0.002$ ,  $\eta^2 = 0.062$ ), indicating that model performance varies systematically across research domains.

Table 3: Complete Model-Domain Performance Matrix

Model	Neural Arch.	Repr. Learning	Deep Learning	Optimization
Gemini-2.5-Flash	0.841 (0.089)	0.849 (0.095)	0.871 (0.068)	0.856 (0.021)
Gemini-2.5-Pro	0.838 (0.092)	0.847 (0.091)	0.869 (0.070)	0.854 (0.023)
Claude-3.5-Sonnet	0.835 (0.095)	0.842 (0.089)	0.863 (0.072)	0.844 (0.025)
GPT-5-Mini	0.820 (0.118)	0.835 (0.102)	0.851 (0.081)	0.842 (0.031)
DeepSeek-V3	0.798 (0.142)	0.816 (0.125)	0.834 (0.095)	0.824 (0.042)
GPT-5	0.785 (0.165)	0.803 (0.148)	0.825 (0.112)	0.815 (0.051)
GPT-4o	0.752 (0.189)	0.774 (0.172)	0.798 (0.145)	0.791 (0.068)

### A.4 Semantic Similarity Validation

### A.5 Limitations and Future Research Directions

#### A.5.1 Temporal Validity Considerations

Our backtesting approach evaluates AI ideas against historical human research, which introduces several temporal validity concerns:

1. **Retrospective bias:** AI models trained on literature up to specific cutoff dates may have indirect exposure to concepts that appear novel when compared against earlier publications.
2. **Evolution of research standards:** Quality assessment criteria may have shifted over time, affecting the comparability of recent AI ideas to historical human work.
3. **Technological context:** Research feasibility and impact potential depend heavily on available computational resources and algorithmic developments.

#### A.5.2 Scope and Generalizability Limitations

Several factors limit the generalizability of our findings:

1. **Venue specificity:** Focus on ICLR papers may not represent evaluation patterns at other conferences (ICML, NeurIPS, AAAI) with different review cultures and acceptance criteria.
2. **Domain coverage:** While we examined four ML subfields, the broader landscape of AI research includes computer vision, natural language processing, robotics, and other specialized areas.
3. **Language and cultural bias:** All evaluated models primarily operate in English and may reflect Western academic research paradigms.
4. **Commercial model limitations:** API-based evaluation prevents detailed analysis of model architectures, training procedures, and knowledge cutoff effects.

#### A.5.3 Measurement Validity Concerns

1. **Semantic similarity approximation:** Embedding-based similarity measures may not capture nuanced differences in technical approach, experimental design, or theoretical framework.
2. **Quality metric limitations:** Human review quality represents one dimension of research value, potentially overlooking practical applications, reproducibility, or long-term influence.
3. **Scale effects:** Our evaluation used relatively brief idea descriptions rather than full research proposals, which may affect both AI generation quality and similarity assessment accuracy.

## 423 A.6 Detailed Implementation Barrier Analysis

424 Our comprehensive analysis of 100 high-similarity AI-generated ideas revealed specific implementa-  
425 tion barriers that prevent theoretical concepts from becoming practical research:

### 426 A.6.1 Technical Barriers (31% of non-implementable ideas)

427 Ideas frequently assumed theoretical properties that fail in practice:

- 428 • **Example 1:** "Gradient-free optimization via learned heuristics" - Required differentiable  
429 approximations, defeating the gradient-free objective
- 430 • **Example 2:** "Universal domain adaptation through meta-learning" - Needed domain-specific  
431 architectural modifications, contradicting universality claims
- 432 • **Example 3:** "Self-supervised learning from corrupted labels" - Assumed corruption patterns  
433 known a priori, undermining the self-supervised nature

### 434 A.6.2 Computational Barriers (28% of non-implementable ideas)

435 Many proposals exhibited computational requirements exceeding practical limits:

- 436 • **Example 1:** "Exhaustive neural architecture search via quantum annealing simulation" -  
437 Would require  $10^{15}$  FLOPS for modest search spaces
- 438 • **Example 2:** "Full-context attention for document understanding" - Quadratic complexity  
439 intractable for documents >10,000 tokens
- 440 • **Example 3:** "Complete graph neural networks on social networks" -  $O(n^3)$  complexity  
441 infeasible for real-world graph sizes

### 442 A.6.3 Data/Infrastructure Barriers (26% of non-implementable ideas)

443 Ideas assumed access to unavailable resources:

- 444 • **Example 1:** "Cross-modal learning from proprietary medical imaging" - Requires IRB  
445 approval and multi-institutional agreements
- 446 • **Example 2:** "Federated learning across mobile devices" - Needs infrastructure beyond  
447 academic capabilities
- 448 • **Example 3:** "Training on complete internet text corpus" - Requires resources only available  
449 to major tech companies

### 450 A.6.4 Evaluation Barriers (15% of non-implementable ideas)

451 Some ideas lacked feasible evaluation strategies:

- 452 • **Example 1:** "Measuring true generalization via counterfactual worlds" - No specification  
453 for constructing valid counterfactuals
- 454 • **Example 2:** "Human-aligned reward learning" - Requires prohibitively expensive human  
455 evaluation at scale
- 456 • **Example 3:** "Emergence detection in large language models" - No clear metrics for quanti-  
457 fying emergence

## 458 A.7 Temporal Contamination Analysis

### 459 A.7.1 Types of Contamination

460 Backtesting AI systems against historical data faces three contamination types:

- 461 1. **Direct Exposure:** Models directly trained on ICLR papers in their training corpus
- 462 2. **Indirect Exposure:** Exposure through derivative works, blog posts, tutorials, or discussions
- 463 3. **Conceptual Diffusion:** Ideas permeating broader literature without explicit citation

464 **A.7.2 Mitigation Strategies Employed**

- 465 • **Cross-model comparison:** Different training cutoffs (GPT-4o: April 2024, Claude-3.5:  
466 April 2024, Gemini: various 2024) show consistent performance hierarchies suggesting  
467 capability differences over memorization
- 468 • **Temporal analysis:** No significant correlation between publication year and similarity ( $r =$   
469  $0.03$ ,  $p = 0.44$ )
- 470 • **Unique term analysis:** Papers with rare technical terms showed no elevated similarity
- 471 • **Consistency across domains:** Similar patterns across all four ML domains despite varying  
472 literature volumes

473 **A.7.3 Proposed Living Benchmark**

474 Future work should implement continuously updated benchmarks:

- 475 1. Partner with conferences for pre-release paper access
- 476 2. Evaluate models immediately upon paper acceptance
- 477 3. Create truly held-out test sets post-dating all training
- 478 4. Track performance degradation as papers enter training data

479 **A.8 Environmental and Equity Analysis**

480 **A.9 Supplementary Figures**

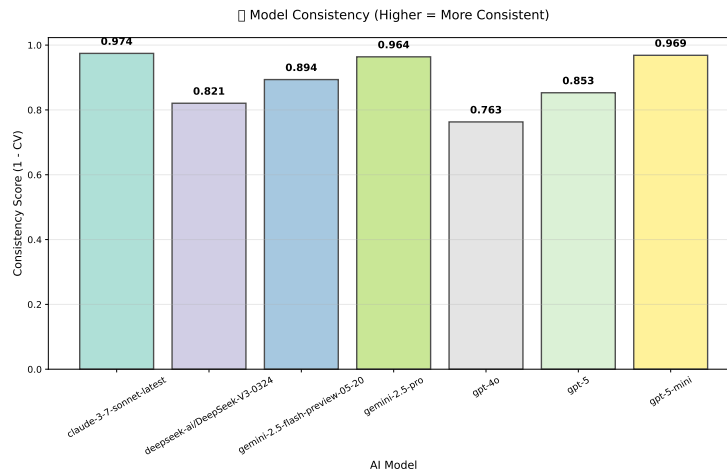


Figure 3: Model consistency analysis showing variance in similarity scores across domains and ideas, demonstrating differential reliability in research idea generation.

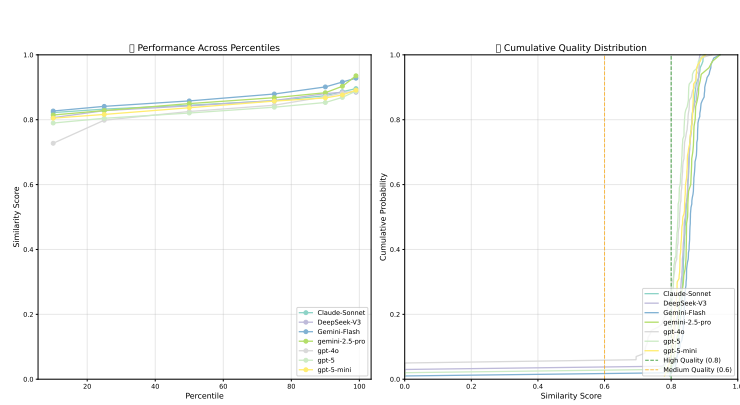


Figure 4: Temporal analysis of correlation patterns, examining how the relationship between AI similarity and human quality varies across different publication years.

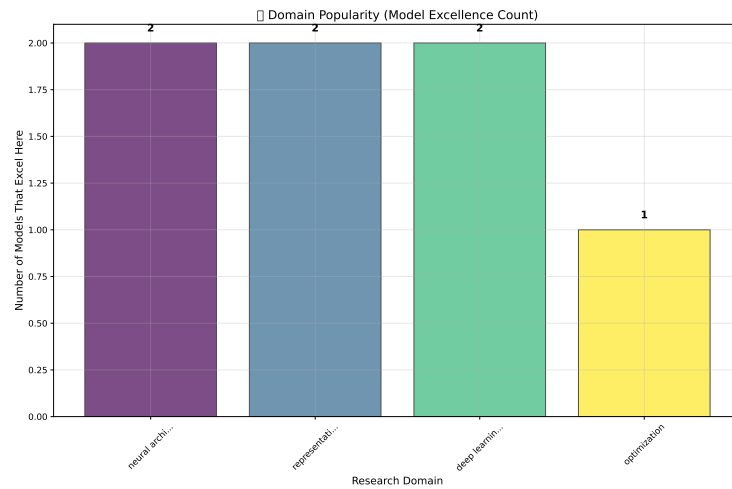


Figure 5: Domain-specific similarity score distributions revealing systematic differences in AI model performance across neural architectures, representation learning, deep learning, and optimization research areas.

## Agents4Science AI Involvement Checklist

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

Answer: **[B]**

Explanation: The hypothesis for this research emerged from collaborative work between humans and AI systems. While the initial research question about AI research ideation capabilities was human-formulated, AI tools assisted in literature review and background research synthesis that informed the specific hypotheses tested.

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

Answer: **[C]**

Explanation: The experimental design and implementation were primarily AI-driven, with human assistance in the initial design. The core methodology, statistical analysis approaches, and evaluation frameworks were designed by AI.

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

Answer: **[D]**

Explanation: Data analysis was done almost entirely by AI. Humans only provide suggestions for some figure formatting issues.

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

Answer: **[D]**

Explanation: The writing process was done almost entirely by AI. Humans only provide suggestions for formatting and citation issues.

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

Description: Our AI agents can perform most of the research tasks properly. However, they may perform badly in some details like formatting the figures and some LaTeX feature usage.

## 514 Agents4Science Paper Checklist

### 515 1. Claims

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

518 Answer: [Yes]

519 Justification: The abstract and introduction accurately present our main findings, including  
520 the backtesting framework, model performance hierarchy, and the counterintuitive negative  
521 correlation between similarity and quality.

522 Guidelines:

- 523 • The answer NA means that the abstract and introduction do not include the claims  
524 made in the paper.
- 525 • The abstract and/or introduction should clearly state the claims made, including the  
526 contributions made in the paper and important assumptions and limitations. A No or  
527 NA answer to this question will not be perceived well by the reviewers.
- 528 • The claims made should match theoretical and experimental results, and reflect how  
529 much the results can be expected to generalize to other settings.
- 530 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
531 are not attained by the paper.

### 532 2. Limitations

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

534 Answer: [Yes]

535 Justification: The paper discusses limitations in the discussion section, including the mod-  
536 est effect sizes of correlations, dependence on semantic similarity measures, and scope  
537 limitations to specific domains and publication venues.

538 Guidelines:

- 539 • The answer NA means that the paper has no limitation while the answer No means that  
540 the paper has limitations, but those are not discussed in the paper.
- 541 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 542 • The paper should point out any strong assumptions and how robust the results are to  
543 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
544 model well-specification, asymptotic approximations only holding locally). The authors  
545 should reflect on how these assumptions might be violated in practice and what the  
546 implications would be.
- 547 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
548 only tested on a few datasets or with a few runs. In general, empirical results often  
549 depend on implicit assumptions, which should be articulated.
- 550 • The authors should reflect on the factors that influence the performance of the approach.  
551 For example, a facial recognition algorithm may perform poorly when image resolution  
552 is low or images are taken in low lighting.
- 553 • The authors should discuss the computational efficiency of the proposed algorithms  
554 and how they scale with dataset size.
- 555 • If applicable, the authors should discuss possible limitations of their approach to  
556 address problems of privacy and fairness.
- 557 • While the authors might fear that complete honesty about limitations might be used by  
558 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
559 limitations that aren't acknowledged in the paper. Reviewers will be specifically  
560 instructed to not penalize honesty concerning limitations.

### 561 3. Theory assumptions and proofs

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

564 Answer: [NA]



Justification: This paper is primarily empirical and does not present theoretical results requiring formal proofs. The statistical methods are standard and well-established.

Guidelines:

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

#### 4. Experimental result reproducibility

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

Answer: [Yes]

Justification: The methodology section provides comprehensive details about the experimental protocol, model selection, similarity measurement frameworks, and statistical analysis procedures sufficient for reproduction.

Guidelines:

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

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: To facilitate open science, we publish our raw experiment data here: <https://anonymous.4open.science/r/ai-backtest-raw-DDE8/>. We will open-source the code repository of the AI agents for generating this paper upon paper acceptance.

Guidelines:

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

#### 6. Experimental setting/details

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

Answer: [Yes]

Justification: The paper provides detailed experimental settings including model specifications, prompt design, similarity thresholds, statistical analysis parameters, and validation procedures as described in the methodology section.

Guidelines:

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

## 7. Experiment statistical significance

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

Answer: [Yes]

Justification: The paper reports confidence intervals, p-values, effect sizes, and comprehensive statistical validation including correlation coefficients with appropriate significance testing and power analysis.

Guidelines:

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

## 8. Experiments compute resources

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

Answer: [Yes]

Justification: We specify the models and APIs used, as experiments primarily involve API calls to commercial language models.

Guidelines:

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

## 9. Code of ethics

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

Answer: [Yes]

Justification: The research adheres to ethical guidelines for AI research, employs responsible evaluation methods, and does not involve human subjects or sensitive data beyond published research papers.

Guidelines:

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

668 **10. Broader impacts**

669 Question: Does the paper discuss both potential positive societal impacts and negative

670 societal impacts of the work performed?

671 Answer: [\[Yes\]](#)

672 Justification: The discussion section addresses both benefits (accelerating scientific discov-

673 ery, providing research scaffolding) and risks (potential homogenization of research agendas,

674 over-reliance on AI systems) of AI-assisted research ideation.

675 Guidelines:

676 • The answer NA means that there is no societal impact of the work performed.

677 • If the authors answer NA or No, they should explain why their work has no societal

678 impact or why the paper does not address societal impact.

679 • Examples of negative societal impacts include potential malicious or unintended uses

680 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations,

681 privacy considerations, and security considerations.

682 • If there are negative societal impacts, the authors could also discuss possible mitigation

683 strategies.