

---

# ScientistOne: Verifiable Autonomous Research via Chain-of-Evidence

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Autonomous research agents produce competitive solutions and complete  
2 manuscripts, yet their outputs suffer from verifiability failures that are largely  
3 invisible to evaluations that check fluency rather than evidence: fabricated citations,  
4 unreproducible scores, and method descriptions that diverge from the implementa-  
5 tion. These failures share a common root: no existing system is designed to trace  
6 claims back to evidence, and no existing evaluation protocol audits whether the  
7 claims are supported. We formalize this gap through *Chain-of-Evidence* (CoE), a  
8 verifiability standard requiring every claim to be traceable to its evidence source,  
9 and instantiate it in two contributions: *ScientistOne*, an end-to-end autonomous  
10 research system that maintains evidence chains by construction throughout litera-  
11 ture review, solution discovery, and paper writing; and *CoE Audit*, an evaluation  
12 protocol whose four integrity checks—score verification, specification violation,  
13 citation verification, and method–code alignment—apply uniformly to all systems.  
14 Across 60 papers from four systems on five frontier systems-research tasks, we  
15 find that every baseline exhibits at least one systematic failure mode: phantom  
16 citation rates reach 21%, score verification passes in as few as 42% of papers,  
17 and method–code alignment ranges from 20% to 80%. *ScientistOne* is the only  
18 system to achieve zero phantom citations (0/337 references), perfect score verifi-  
19 cation (12/12), and the highest method–code alignment (14/15), while remaining  
20 competitive on solver performance. We further demonstrate that *ScientistOne*  
21 generalizes to six additional tasks spanning medical imaging, fine-grained recog-  
22 nition, 3D perception, and parameter-constrained language modeling, achieving  
23 state-of-the-art on Parameter Golf and gold medals on MLE-Bench tasks where  
24 baselines fail entirely.

## 25 1 Introduction

26 Large language models are increasingly deployed not as isolated assistants but as autonomous agents  
27 that conduct entire research workflows—from literature review and hypothesis generation through  
28 experimental design and execution to manuscript writing [1, 2, 3, 4, 5, 6, 7]. On systems-optimization  
29 tasks, such agents now produce solutions competitive with human experts [8, 9], and end-to-end  
30 pipelines have generated papers accepted at peer-reviewed workshops [2]. The resulting artifacts—  
31 code, experimental results, and complete manuscripts—are increasingly difficult to distinguish from  
32 human-authored research on surface quality alone.

33 This rapid capability growth exposes a structural tension between *generation* and *verification*. Au-  
34 tonomous research systems operate as multi-stage pipelines in which each stage consumes the output  
35 of the previous one: a literature summary shapes the hypothesis, the hypothesis determines the exper-  
36 iment, and experimental results feed into the manuscript. In such architectures, errors introduced at  
37 any stage are not merely preserved but amplified—a flawed summary can bias experimental design,

38 and a misinterpreted result can propagate into a paper that appears internally coherent precisely  
39 because the same error threads through every section. The risk grows with trajectory length: agents  
40 struggle to track an ever-expanding context [10, 11], hallucinate, and lose sight of the original purpose.  
41 The problem is compounded by fundamental limitations in how language models handle evidence:  
42 generated text is difficult to verify against sources [12], factual claims drift from their grounding [13],  
43 and scientific citations are frequently inaccurate or fabricated [14].

44 In autonomous pipelines, these failure modes interact and compound—a model can embellish  
45 method descriptions beyond what the code implements, report scores that do not reproduce under  
46 the benchmark’s own evaluator, and populate bibliographies from parametric memory rather than  
47 retrieval, all while producing prose that reads as technically sound. Existing evaluation protocols,  
48 whether automated review scores or benchmark leaderboards, reward output fluency and procedural  
49 completion but do not assess whether individual claims trace to supporting evidence.

50 This verifiability gap is not hypothetical. In a systematic audit of 60 papers generated by four  
51 autonomous research systems across five benchmark tasks, we find that *every baseline paper* exhibits  
52 at least one evidence chain failure: phantom citations that do not correspond to any real publication  
53 (up to 21% of all references), method sections that describe algorithms entirely absent from the  
54 submitted code, unreproducible scores, and solution code that exploits the evaluator rather than  
55 solving the task. These failures share a common root cause: *no existing system is designed to trace*  
56 *claims back to evidence, and no existing evaluation protocol audits whether the claims are supported.*

57 We address this with *Chain-of-Evidence* (CoE), a verifiability standard for AI-driven research. Just  
58 as ACID<sup>1</sup> [15] defines what “reliable” means as a standard for a database transaction, CoE defines  
59 what “verifiable” means for a research claim: **every claim must trace, through a recorded evidence**  
60 **chain, to a grounding source.** We instantiate CoE in three ways:

- 61 1. **The CoE Standard (§3)**: a claim taxonomy (citation, numerical, methodological, conclusion)  
62 and the evidence chain structure required for each type.
- 63 2. **ScientistOne (§4)**: an end-to-end autonomous research system whose pipeline—Problem  
64 Investigator, Discovery Engine, and Paper Writer with Claim Verifier—is designed to satisfy CoE  
65 natively. The Problem Investigator reads up to 100 full-text PDFs per topic, producing grounded  
66 experiment briefs; the Claim Verifier checks every claim in the draft against its declared evidence  
67 source before the final paper is produced.
- 68 3. **CoE Audit (§5)**: a post-hoc evaluation protocol for auditing an AI-driven research paper through  
69 four integrity checks—Score Verification, Specification Violation, Citation Verification, and  
70 Method-Code Alignment—targeting the most damaging evidence chain failures.

71 We apply CoE Audit to 15 papers from each of four systems across five frontier systems-research tasks  
72 from ADRS [8, 16] (§6). Every baseline exhibits at least one integrity check failure. *ScientistOne*  
73 achieves zero phantom citations (0/337 references), perfect score verification (12/12), and the highest  
74 method–code alignment (14/15), while remaining competitive on solver performance across all  
75 five tasks. We further demonstrate that *ScientistOne* generalizes to six additional tasks spanning  
76 medical imaging, fine-grained recognition, 3D perception, and parameter-constrained language  
77 modeling, achieving state-of-the-art on Parameter Golf and gold medals on MLE-Bench tasks where  
78 baselines fail entirely.

## 79 2 Related Work

80 **Autonomous research agents.** End-to-end autonomous research systems have rapidly expanded  
81 from constrained ML templates to multi-stage pipelines that coordinate literature grounding, hy-  
82 pothesis generation, experimentation, and paper writing. The AI Scientist [1] pioneered end-to-end  
83 automation but operates on fixed ML templates with hallucinated references and workshop-reject  
84 quality. AI Scientist-v2 [2] advances this with agentic tree search over experimental branches and  
85 review-aware reporting, achieving stronger results on NeurIPS-style tasks. Several concurrent sys-  
86 tems extend the pipeline in different directions: PiFlow [6] introduces principle-aware workflow  
87 coordination; Curie [17] emphasizes execution-intensive experimentation with controlled reruns and  
88 ablations; CodeScientist [7] grounds ideation in open-source code repositories; Agent Laboratory [5]  
89 provides a template-bounded full-loop environment; and AlphaEvolve [9] and EvoScientist [18] apply

---

<sup>1</sup>Atomicity, consistency, isolation, durability.

90 evolutionary search to scientific optimization. Despite this architectural diversity, a common pattern  
91 emerges: generation and execution capabilities have scaled faster than validation and provenance  
92 mechanisms, so systems that produce polished manuscripts may still lack credible scientific closure.  
93 `ScientistOne` targets this gap—rather than advancing the autonomy frontier, we focus on making  
94 autonomous research outputs verifiable.

95 **LLM-driven optimization and benchmarks.** The ADRS benchmark [8] collects real frontier  
96 computer system research questions and serves as our primary evaluation testbed. Evolutionary  
97 search frameworks (AlphaEvolve [9], AdaEvolve [19], EvoScientist [18]) achieve strong results  
98 on ADRS by only focusing algorithm discovery and implementation optimization, while not per-  
99 forming literature grounding, novelty assessment, nor paper writing. Broader evaluation resources  
100 have recently proliferated: Auto-Bench [20], ResearchBench [21], and ResearcherBench [22] probe  
101 discovery-oriented research behavior; MAgentBench [23], EXP-Bench [24], and PaperBench [25]  
102 stress-test experimentation, replication, and execution reliability; while AIRS-Bench [26] and FIRE-  
103 Bench [27] target frontier research tasks and rediscovery verification. However, most benchmarks  
104 measure procedural correctness—whether a system can complete a research-like workflow—rather  
105 than whether the resulting claims are actually supported by evidence. CoE Audit fills this gap by  
106 operationalizing validity, reliability, and provenance as a cross-system audit protocol. We also show-  
107 case `ScientistOne` that generalizes beyond computer system research problems to other domains,  
108 including computer vision, healthcare [28] and low-bit quantization on language modeling [29].

109 **Scientific integrity and provenance.** Current autonomous research systems range in reporting  
110 maturity from direct manuscript drafting [1, 7, 3] to iterative review-aware revision [2, 30] to artifact-  
111 linked reporting where manuscripts are tied to execution traces and code [4, 17]. Artifact-linked  
112 reporting provides the strongest traceability but remains uncommon; most systems produce fluent  
113 papers that mask broken evidence chains.

114 Prior work on citation hallucination [12], factual accuracy [13], and citation-level claim support [14]  
115 performs post-hoc detection at the text level. CoE differs in two ways: it defines a *constructive*  
116 standard (evidence chains must be built at claim-production time, not recovered after the fact),  
117 and it covers the full research artifact (paper + code + evaluator logs), not just text. CoE Audit  
118 operationalizes this standard as a cross-system evaluation protocol that can be applied in forensic  
119 mode to any system’s released artifacts.

### 120 3 Chain-of-Evidence: A Verifiability Standard

121 ***Principle:** Every claim produced by a research system must be traceable, through*  
122 *a chain of recorded evidence, to a grounding source—and that chain must be cheap*  
123 *to walk for both humans and machines.*

124 A credible research claim must be backed by verifiable evidence. Without this requirement, the same  
125 system that produces a plausible-sounding paper can also produce fabricated citations, hallucinated  
126 numbers, and descriptions of experiments that never happened. Just as a database that violates ACID  
127 may return plausible-looking query results even as it silently corrupts data, a research system that  
128 violates CoE may produce plausible-looking papers whose claims cannot be traced to evidence. ACID  
129 does not prescribe *how* to build a database; it prescribes what properties the database must have. CoE  
130 similarly prescribes what properties a research artifact must have, independent of system architecture.

131 We define four claim types, each with a required evidence chain shape. **Citation claims** (e.g., “Smith  
132 et al. showed X”) require that the cited work exists in a scholarly database and that its content is  
133 consistent with the assertion made about it. **Numerical claims** (e.g., “achieves 87.3% on Prism”)   
134 must trace from the reported value to an execution log produced by an evaluator run. **Methodological**  
135 **claims** (e.g., “we use a 3-layer MLP”) must resolve from the method description to corresponding  
136 source code. **Conclusion claims** (e.g., “outperforms baseline by 5%”) derive from supporting  
137 numerical claims through a verifiable comparison. CoE is deliberately architecture-agnostic: it  
138 defines what properties a verifiable artifact should have, not how the system should construct one.

139 In the following sections, we describe `ScientistOne`, an autonomous research system designed to  
140 satisfy CoE by construction (§4), and CoE Audit, an automated audit framework that measures how  
141 well any system’s artifacts meet the standard through four integrity checks (§5).

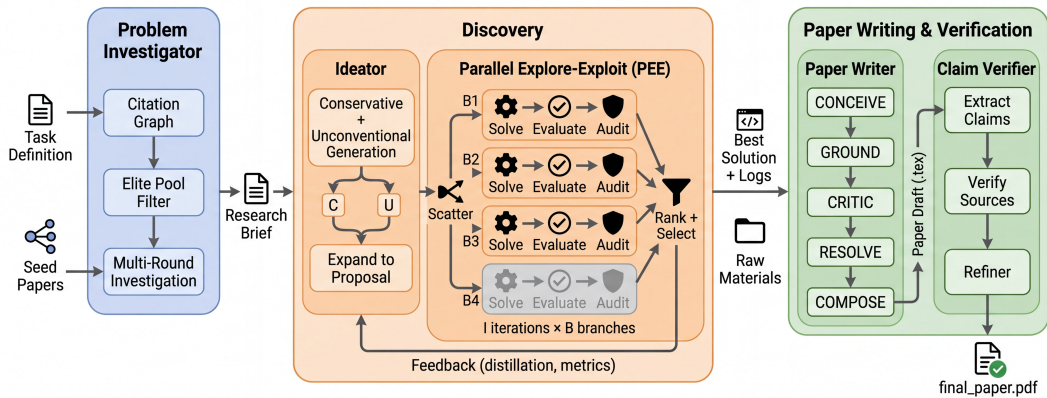


Figure 1: ScientistOne pipeline. Stage 1 grounds the literature via retrieved PDFs; Stage 2 explores and evaluates solutions across parallel branches; Stage 3 writes and verifies the paper, with a Claim Verifier that checks every claim against its evidence source before the final output is produced.

## 142 4 ScientistOne: Research with Verifiability

143 We now describe ScientistOne, an end-to-end autonomous research system whose three-stage  
 144 architecture is shaped by the CoE requirements: each module is designed to produce structured  
 145 artifacts that carry the provenance metadata needed to verify claims against their evidence (Figure 1).

### 146 4.1 Stage 1: Literature Grounding

147 The Problem Investigator (PI) is designed to ensure that every paper the system cites was retrieved  
 148 from a scholarly database, read in full text, and recorded with provenance metadata. Without  
 149 structured retrieval, autonomous systems generate citations from model memory—our audit finds this  
 150 produces phantom references at rates of up to 21% (§6.1). PI addresses this by construction: starting  
 151 from seed papers, it builds a citation graph via the Semantic Scholar API, reads up to 100 full-text  
 152 PDFs per topic, and produces a structured research brief. The brief feeds the Ideator; PI’s seed  
 153 reference bibliography provides grounding material for citation claims in the final paper. Pipeline  
 154 details are in §C.

### 155 4.2 Stage 2: Discovery

156 The Ideator generates candidate approaches from the PI brief, scores them on novelty and feasibility,  
 157 and scatters the top-ranked proposals across parallel branches of the *Parallel Explore-Exploit* (PEE)  
 158 orchestrator. Each branch runs an isolated cycle: a Solver agent implements a solution, a task-specific  
 159 evaluator scores it, and an auditor checks idea–implementation alignment. High-performing branches  
 160 are refined; underperformers are replaced by new explorations. After  $I$  iterations across  $B$  branches,  
 161 all execution logs and evaluator outputs are assembled into `verified_scores.json`—the ground  
 162 truth that the Claim Verifier checks numerical claims against. Architecture details are in §C.

### 163 4.3 Stage 3: Paper Writing & Verification

164 **Paper Writer.** The Paper Writer first builds a structured representation from all raw materials (PI  
 165 brief, experimental log, verified scores, code), where every factual claim carries an inline evidence  
 166 tag binding it to a workspace artifact. Deterministic checks validate score presence and baseline  
 167 attribution; an LLM critic audits story-level coherence; and a resolve loop refines until convergence.  
 168 Per-section writers then render the representation into  $\text{\LaTeX}$ , with each claim committing to its  
 169 evidence source at writing time—verified numbers and named baselines are the first context the writer  
 170 sees, reducing the incentive to fill gaps from model memory. The composed draft passes through the  
 171 Claim Verifier; unsupported claims are corrected or dropped before the final paper is produced.

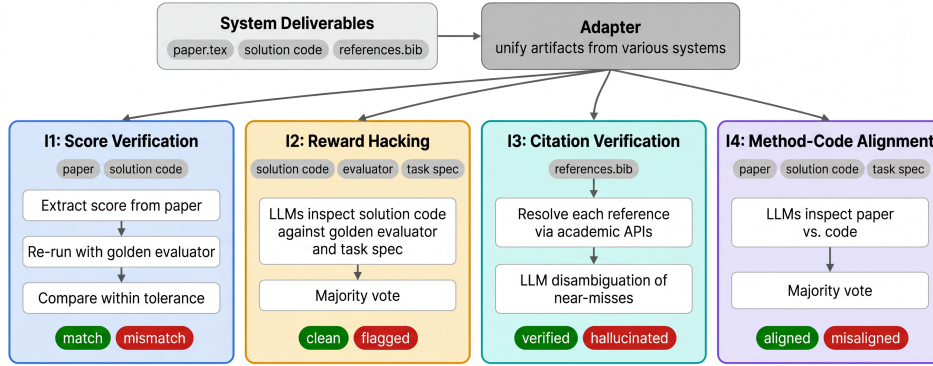


Figure 2: CoE Audit overview.

172 **Claim Verifier.** The Claim Verifier checks every claim in the draft against its declared evidence  
 173 source: numerical claims are compared by tolerance against evaluator logs, citation claims are  
 174 resolved against the bibliography, and methodological claims are checked against experimental logs.  
 175 Unsourced claims are flagged automatically. Per-type verification rules are in §C.3.

## 176 5 CoE Audit: Measuring Verifiability

177 To evaluate verifiability of research outcomes—both ours and baselines’—we introduce *CoE Audit*,  
 178 an evaluation protocol applied identically to any system (Figure 2). An adapter unifies each system’s  
 179 deliverables (paper, solution code, references) into a common artifact bundle, then four integrity  
 180 checks run independently each targeting a specific way evidence chains break:

181 **Score Verification (I1).** The paper’s reported score is extracted by LLMs from both  $\text{\TeX}$  and PDF  
 182 files, then compared against reproduced scores with solution files on the golden evaluator. A paper  
 183 passes if the two match within an adaptive tolerance that accounts for evaluator noise.

184 **Specification Violation (I2).** Specification violations occur when solution code breaks task rules—  
 185 for example, reading evaluator source files to extract expected outputs, or hardcoding answers  
 186 for known test cases. The generating agent optimizes for the score, not for the problem the task  
 187 is designed to measure. LLMs inspect the solution code against the golden evaluator and task  
 188 specification to detect such violations, with majority vote across multiple runs.

189 **Citation Verification (I3).** Each reference is resolved by querying multiple academic APIs (Se-  
 190 mantic Scholar, arXiv, OpenAlex, CrossRef) using arXiv ID, DOI, and title. An LLM cross-checks  
 191 the full bib entry against returned records to catch near-misses and citation gaming (e.g., a real DOI  
 192 attached to a fabricated description). References matching no record are classified as hallucinated.

193 **Method-Code Alignment (I4).** An LLM reads the paper’s method section and the solution code  
 194 side by side, then judges whether the paper faithfully describes what the code does. Acceptable  
 195 simplification (e.g., omitting implementation details) is treated as aligned; only cases where the paper  
 196 describes a fundamentally different algorithm count as misaligned. We conduct multiple independent  
 197 runs with majority vote to reduce LLM judgment noise.

## 198 6 Experiments

199 **Benchmark.** We evaluate on the Automated Design of Research Systems (ADRS) benchmark [8],  
 200 which collect five frontier computer system research problems: **Prism** (GPU memory allocation for  
 201 LLM serving), **Cloudcast** (cloud network cost optimization), **EPLB** (expert-parallel load balancing  
 202 for MoE models), **LLM-SQL** (text-to-SQL query optimization), and **TXN** (transaction scheduling).  
 203 Each task provides a fixed evaluator, starter code, and scoring metric. We choose ADRS as our  
 204 primary benchmark for three reasons: (1) the tasks are drawn from real-world systems-optimization

205 problems and researchers published results on those in the past years, not synthetic puzzles; (2) the  
 206 leaderboard provides both human expert and recent LLM-agent baselines [8], enabling apples-to-  
 207 apples comparison; and (3) the gold-standard evaluators are deterministic enough to support Score  
 208 Verification (§5) and Specification Violation detection (§5). Following [19], who observe that ADRS  
 209 evaluations involve “noisy feedback,” we run each evaluator five times and compare against adaptive  
 210 tolerance  $\max(1\%, 3\sigma/|\bar{s}|)$  to account for inherent evaluator variance.

211 **Baseline Systems.** We evaluate three baseline systems and `ScientistOne`. All baselines are  
 212 open-source, enabling us to adapt each to the ADRS benchmark for controlled comparison under  
 213 identical conditions. They span the design spectrum from highly structured scaffolding to fully  
 214 autonomous agents, providing coverage across the major paradigms for AI research agents.

- 215 • **AutoResearchClaw (ARC)** [31]: 23-stage waterfall pipeline with OpenAlex/Semantic Scholar.
- 216 • **DeepScientist (DS)** [3]: Autonomous research studio with Codex CLI and MCP tools.
- 217 • **AI-Researcher (AIR)** [4]: MetaChain ReAct agents, adapted for systems-optimization.
- 218 • **ScientistOne**: Full pipeline with evidence chain maintenance from problem framing through  
 219 paper composition (§4).

220 We also ran Sakana AI-Scientist v2 [2] on all five ADRS tasks but exclude it from results: its  
 221 pipeline’s hardcoded ML-training stage design (e.g., “introduce datasets from HuggingFace,” “change  
 222 learning rate and batch size”) overrides the task specification, causing agents to train neural networks  
 223 rather than optimize the target functions. Only one of fifteen runs produced a solution that called the  
 224 ADRS evaluator as specified.

225 **ADRS Adaptation.** Each baseline was adapted to the ADRS benchmark with varying levels of  
 226 source modification, from prompt-only changes (DS) to patching 8–19 source files (ARC, AIR)  
 227 to interface with ADRS task specifications, evaluators, and the NeurIPS 2026 paper template. We  
 228 standardized on Gemini 3.1 Pro as the backbone LLM across all systems for both solver code  
 229 generation and paper writing. To ensure best-effort performance, we configured generous iteration  
 230 and timeout budgets: up to 20 solver iterations per task— $6.7\times$  the default for ARC, though most  
 231 tasks converge well before this cap—and 2-hour code generation windows. Runs that crashed due to  
 232 infrastructure issues (API timeouts, rate limits, LaTeX compilation errors) were re-attempted with  
 233 fresh state, up to 3 attempts per run; no run was re-attempted to improve solver scores. Of the 60 total  
 234 runs, 14 required at least one retry. All 60 runs ultimately produced solver code and a compiled paper,  
 235 though artifact quality varies. For each system, we run 3 seeds per task, producing 15 papers per  
 236 system and 60 papers total. CoE Audit is applied identically to all systems. Full adaptation details  
 237 are provided in Section G.

## 238 6.1 CoE Audit Results

Table 1: CoE Audit results across four systems (15 papers per system). EPLB papers are excluded from Score Verif. (no system reports a single Combined Score). Metric definitions are in §5.

System	Score Verif. $\uparrow$	Spec. Violation $\downarrow$	Citation Verif. $\downarrow$	Method-Code $\uparrow$
AutoResearchClaw [31]	5/12	0/15	3/196	3/15
DeepScientist [3]	11/12	1/15	42/201	5/15
AI-Researcher [4]	9/12	1/15	21/222	12/15
ScientistOne	12/12	0/15	0/337	14/15

239 The results of CoE Audit across four systems are presented in Table 1. All I1–I3 flagged positives were  
 240 manually verified and corrected by human reviewers; I4 judgments were validated on a sampled basis.  
 241 `ScientistOne` is the only system that achieves perfect score verification (12/12), zero specification  
 242 violations (0/15), zero phantom citations, and 14/15 method-code alignment. The gap is largest in  
 243 citation integrity and method-code alignment—the two checks that test evidence provenance rather  
 244 than score reproduction.

245 **Score verification and specification violation.** `ScientistOne` achieves perfect score verification  
 246 (12/12): every paper’s claimed result reproduces exactly under re-evaluation. DS matches in 11/12

247 (92%), with the single failure caused by fabricated metric direction—the paper relabels a cost  
 248 metric as “utility” and inverts the optimization direction, so a worse-than-baseline score reads as  
 249 an improvement; rerunning the declared solution against the true metric confirms the baseline still  
 250 wins. AIR matches in 9/12 (75%), with the failures being small-magnitude discrepancies (1–4%)  
 251 plus one paper that reports no quantitative scores at all. ARC matches in 5/12 (42%), and its failures  
 252 span a wider range—from a few percent up to a  $2\times$  gap—and include qualitatively distinct cases: a  
 253 paper without any results table, submission code that prevents clean re-evaluation, and a paper that  
 254 acknowledges an integration failure with the benchmark evaluator. Specification violation rates are  
 255 uniformly low: ARC and `ScientistOne` show 0/15, while AIR and DS each have one flagged paper  
 256 involving evaluator file modifications detected during code inspection (see Case 3, §B.1).

257 **Citation integrity.** `ScientistOne` achieves zero phantom citations across all 15 papers (0 out of  
 258 337 unique references). DS exhibits the highest phantom rate (42/201, 20.9%), followed by AIR  
 259 (21/222, 9.5%) and ARC (3/196, 1.5%). ARC’s low phantom rate reflects its structured retrieval  
 260 pipeline; its three phantoms are a single repeated malformed bib entry across three papers. DS and  
 261 AIR rely on model memory for reference generation, producing plausible-looking but non-existent  
 262 citations—a failure mode illustrated in Case 2 (§B.1). `ScientistOne`’s zero phantom rate is an  
 263 architectural property of the PI’s citation graph: every reference originates from a Semantic Scholar  
 264 API call whose result is cached in the evidence chain. (Full list of phantom citations in § E.2.)

265 **Method-code alignment.** `ScientistOne` achieves 14/15 aligned papers (93%) under the MCA-  
 266 0/MCA-1 threshold, compared to AIR (12/15, 80%), DS (5/15, 33%), and ARC (3/15, 20%). The  
 267 misalignment patterns differ qualitatively across systems. AIR’s failures are typically algorithm  
 268 mismatches—e.g., the paper describes a more sophisticated procedure than the code implements.  
 269 ARC exhibits the broadest failure-mode mix: algorithm-class mismatches (e.g. beam search with  
 270 Edmonds’ arborescence vs. greedy edge penalization), undisclosed environment-variable switches  
 271 that swap in entirely different solvers at runtime, evaluation hacking (e.g. returning empty column  
 272 orderings to maximize a prefix-hit metric), and method-vs-ablation inversion between paper and  
 273 final solution code. `ScientistOne`’s single misaligned paper is a case where the paper writer  
 274 fabricated algorithmic claims not present in the code—describing a “hybrid neuro-symbolic solver”  
 275 with “LLM-guided evolutionary search” when the submitted code is a deterministic routing heuristic  
 276 with no LLM calls. The Claim Verifier’s method-code cross-check (§4.3) catches nearly all such  
 277 misrepresentation before paper finalization, though complete coverage remains a limitation (§A).

278 We present detailed Failure Mode case studies in § B.1. We also present sub-categorization and  
 279 statistics for CoE audit errors in § E.1, E.2, E.3 for I1, I3 and I4 respectively.

## 280 6.2 Review Scores

281 We evaluate perceived paper quality using ScholarPeer [32], an automated peer review system backed  
 282 by `gemini-3.1-pro-preview` and rich literature search support.

Table 2: ScholarPeer review rating scores (1–4 scales except Overall (1-10 scale)) and accept decisions. *Average*: mean across 15 papers per system. *Best-of-3*: strongest seed per task.

System	Soundness	Originality	Quality	Clarity	Overall	#Accept
<i>Average (3 seeds <math>\times</math> 5 tasks)</i>						
ARC	1.1	2.3	1.1	2.5	1.9	0/15
DS	1.7	1.7	1.6	3.1	2.5	1/15
AIR	1.9	2.4	1.9	3.1	3.4	2/15
<code>ScientistOne</code>	<b>2.3</b>	<b>2.5</b>	<b>2.3</b>	3.0	<b>4.5</b>	<b>6/15</b>
<i>Best-of-3 (strongest seed per task)</i>						
ARC	1.2	2.0	1.2	3.0	3.0	0/5
DS	2.2	2.0	2.2	3.6	3.6	1/5
AIR	2.0	2.4	2.0	3.2	4.0	1/5
<code>ScientistOne</code>	<b>2.8</b>	<b>3.0</b>	<b>2.8</b>	<b>3.6</b>	<b>6.6</b>	<b>4/5</b>

283 **Verifiable papers are deemed better by automatic reviewers.** `ScientistOne` achieves a 40%  
 284 accept rate (6/15), tripling the best baseline (AIR: 13%), and best-of-3 selection reaches 6.6 overall

Table 3: Solution discovery performance on ADRS benchmark tasks (best-of-3 seeds). ARC/AIR/DS/ScientistOne scores are from independent gold evaluator re-runs on submitted solver code; Human, AdaEvolve, and EvoX scores are from original publications. \* indicates Gemini-3.0-Pro; all other systems use Gemini-3.1-Pro.

Task	Dir.	Human	AdaEvo*	EvoX*	ARC	AIR	DS	ScientistOne	ScientistOne*
Prism	↑	21.89	<b>26.26</b>	<b>26.26</b>	26.25	<b>26.26</b>	<b>26.26</b>	<b>26.26</b>	<b>26.26</b>
Cloudcast	↓	626.24	637.10	623.69	690.37	734.28	620.09	<b>618.08</b>	<b>618.08</b>
EPLB	↑	0.1265	0.1450	0.1453	0.1266	0.1449	0.1284	0.1459	<b>0.1461</b>
LLM-SQL	↑	0.6920	<b>0.7520</b>	0.7300	0.6757	0.7148	0.7307	0.7222	0.7115
TXN	↑	2724.8	4310	4310	3247	<b>4311</b>	4286	3906	3861

285 rating score and 4/5 tasks accepted. This gap is not driven by better algorithms—solver scores cluster  
 286 tightly across systems (Table 3)—but by what happens *after* the solver finishes. The Claim Verifier  
 287 prevents the most damaging failure mode we observe in rejected papers: claims that contradict the  
 288 paper’s own data (e.g., “sub-millisecond latency” when the results table reports 7.9 ms). Notably,  
 289 CoE Audit and review scores measure orthogonal properties: across all 60 papers, Clarity averages  
 290 1–2 points above Soundness, confirming that narrative coherence is independent of evidence integrity.

291 **The paper quality is bottlenecked by research soundness, not writing capability.** Across all  
 292 systems, Clarity is consistently the highest-scoring dimension (2.5–3.1) while Soundness is the lowest  
 293 (1.1–2.3): these papers read well but do not withstand methodological scrutiny. The reviewer’s  
 294 two most frequent complaints are missing comparisons against published baselines and proxy-only  
 295 evaluation without end-to-end system measurements. While ScientistOne’s Problem Investigator  
 296 retrieves related work and identifies candidate baselines, the resulting comparisons do not yet meet  
 297 the depth that ScholarPeer expects (e.g., re-implementing a SOTA method and reporting head-to-head  
 298 numbers). ScientistOne also exhibits high seed variance (e.g., EPLB scores 1, 3, 8 across three  
 299 seeds on the same task): rejected runs are those where the paper writer generates claims unsupported  
 300 by the results (e.g., “sub-millisecond” when latency is 7.9 ms), while accepted runs make calibrated  
 301 claims from the same underlying data. A post-hoc paper review that checks every quantitative claim in  
 302 the abstract against the results tables would catch these contradictions and likely reduce this variance.

### 303 6.3 Solution Discovery Performance

304 To compare the effectiveness of our discovery module with baseline systems, we report ADRS  
 305 performance in Table 3. Following [19], we report best-of-3 scores across seeds for each system.  
 306 For ARC, AIR, DS, and ScientistOne, each score is from independent gold evaluator re-runs on  
 307 submitted solver code, ensuring cross-system comparability. Human, AdaEvolve, and EvoX scores  
 308 (Gemini-3.0-Pro, best-of-3 runs) are from original publications [19, 33]; we did not have access to  
 309 their solver code for independent verification.

310 All six systems substantially outperform the human baseline on every task and fall within a narrow  
 311 performance range, consistent with the observation of [8] that LLM-based agents rapidly converge  
 312 to similar solution quality. Critically, verifiability does not sacrifice performance: ScientistOne  
 313 achieves the best score on Cloudcast and EPLB, and remains competitive on the other three tasks.

## 314 7 Generalizability: MLE-Bench and Parameter Golf

315 To test whether the discovery loop transfers beyond ADRS, we evaluate ScientistOne unmodified  
 316 across a diverse set of six tasks selected for their rigorous demands and relevance to current AI  
 317 research. First, we select five MLE-Bench [28] Kaggle competitions spanning medical imaging,  
 318 fine-grained recognition, and 3D perception. We specifically targeted tasks in the Medium and  
 319 High difficulty tiers because they possess sufficient complexity and substantive research value to  
 320 serve as meaningful benchmarks for automated paper writing and scientific discovery. Second, to  
 321 assess adaptability in a novel, live environment, we evaluate ScientistOne on the Parameter Golf  
 322 competition [29]. This live competition serves as an ideal AI research task focused on LLM training,  
 323 requiring participants to train the highest-performing language model under multi-faceted constraints.  
 324 Full task details and evaluation setups are provided in §F.

Table 4: Comparison of solver performance across five MLE-Bench tasks and Parameter Golf. “Above/Below Median” denotes outperforming or underperforming the median participant. MLE-Bench medals (Gold, Silver, Bronze) represent simulated private leaderboard standings. “SOTA” indicates top-1 performance on the Parameter Golf Leaderboard (as of 2026-04-27).

Task	Dir.	DeepScientist		ScientistOne	
		Score	Highlight	Score	Highlight
3D Object Detection	↑	0.0000	Below Median	0.1763	Gold Medal
AI4Code	↑	0.6964	Below Median	0.8356	Above Median
iMet 2020 FGVC7	↑	0.6804	Silver Medal	0.6791	Silver Medal
RSNA Brain Tumor	↑	0.6377	Gold Medal	0.6518	Gold Medal
iNaturalist 2019 FGVC6	↓	0.2158	Silver Medal	0.2445	Silver Medal
Parameter Golf	↓	Invalid	Size limit exceeded	1.0600	SOTA (Constraints met)

325 **Solver performance generalizes and demonstrates robustness.** As shown in Table 4,  
326 `ScientistOne` demonstrates strong generalization and remarkable robustness across diverse  
327 domains, difficulty levels, and strict constraints. On the High difficulty MLE-Bench tasks,  
328 `ScientistOne` earns two Gold Medals (RSNA Brain Tumor and 3D Object Detection), notably solv-  
329 ing the 3D Object Detection task with a Gold Medal score while `DeepScientist` failed entirely (scoring  
330 0.0000). On the Medium difficulty tasks, `ScientistOne` secures Silver Medals on both iMet 2020  
331 and iNaturalist 2019—remaining highly competitive with `DeepScientist`—and achieves an Above  
332 Median standing on AI4Code, marking a significant improvement. Furthermore, `ScientistOne`  
333 proves its robustness on the completely new LLM training task of Parameter Golf. While the baseline  
334 fails to produce a valid submission due to exceeding the 16MB artifact size limit, `ScientistOne`  
335 adheres to all constraints and achieves state-of-the-art performance with a score of 1.0600.

336 **Solution novelty comparison on Parameter Golf.** Both `ScientistOne` and `DeepScientist` were  
337 provided with the same prior-art reference and achieved superficially similar numerical improvements.  
338 However, the systems arrived at these outcomes in completely different ways. `ScientistOne` demon-  
339 strated genuine research capability by introducing novel algorithmic techniques to the quantization  
340 block: specifically, a Hessian-diagonal-weighted SVD initialization and an alternating-least-squares  
341 (ALS) refinement loop that utilizes GPTQ and a Cholesky-weighted truncated SVD. These mech-  
342 anisms are entirely novel to the leaderboard, and internal ablations isolate the ALS loop as the  
343 primary driver of the performance gain. In contrast, `DeepScientist` contributed zero algorithmic  
344 novelty. Its modifications were limited to superficial environment and portability adjustments rather  
345 than structural research advancements. As a result, `DeepScientist` failed to improve the underlying  
346 algorithm – merely replicating the reference’s performance – and ultimately produced an invalid  
347 submission by exceeding the strict 16MB artifact size limit. Full details are available in our generated  
348 paper provided in the supplementary materials.

## 349 8 Conclusion

350 Autonomous research systems have reached the point where solver quality alone no longer differ-  
351 entiates them—multiple systems achieve competitive scores on the same benchmarks with vastly  
352 different approaches. What separates their outputs is whether the resulting paper can be trusted.  
353 Our 60-paper audit shows that no baseline produces papers free of evidence chain failures, and the  
354 failures—phantom citations up to 21%, fictional method sections, scores on the wrong scale—are  
355 invisible to evaluations that assess narrative quality alone.

356 Chain-of-Evidence reframes verifiability as a first-class design constraint. `ScientistOne` demon-  
357 strates that an end-to-end pipeline can maintain evidence chains without sacrificing solver compet-  
358 itiveness, and CoE Audit provides a reusable protocol for auditing any system’s output. The gap  
359 between `ScientistOne`’s results and the baselines confirms that verifiability is architectural: systems  
360 that build evidence chains at claim-production time outperform those that reconstruct grounding after  
361 the fact. The harder problems—verifying citation support, checking conclusion claims, extending to  
362 domains without deterministic evaluators—remain open, but they are well-defined, and addressing  
363 them matters as autonomous research scales.

364 **References**

- 365 [1] Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The  
366 AI scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint*  
367 *arXiv:2408.06292*, 2024.
- 368 [2] Yutaro Yamada, Robert Tjarko Lange, Cong Lu, Shengran Hu, Chris Lu, Jakob Foerster, Jeff  
369 Clune, and David Ha. The AI scientist-v2: Workshop-level automated scientific discovery via  
370 agentic tree search. *arXiv preprint arXiv:2504.08066*, 2025.
- 371 [3] Yixuan Weng, Minjun Zhu, Qiuji Xie, Qiyao Sun, Zhen Lin, Sifan Liu, and Yue Zhang.  
372 Deepscientist: Advancing frontier-pushing scientific findings progressively. *arXiv preprint*  
373 *arXiv:2509.26603*, 2025.
- 374 [4] Jiabin Tang, Lianghao Xia, Zhonghang Li, and Chao Huang. AI-Researcher: Autonomous  
375 scientific innovation. *arXiv preprint arXiv:2505.18705*, 2025.
- 376 [5] Samuel Schmidgall, Yusheng Su, Ze Wang, Ximeng Sun, Jialian Wu, Xiaodong Yu, Jiang Liu,  
377 Michael Moor, Zicheng Liu, and Emad Barsoum. Agent laboratory: Using LLM agents as  
378 research assistants. *Findings of the Association for Computational Linguistics: EMNLP 2025*,  
379 pages 5977–6043, 2025.
- 380 [6] Yingming Pu, Tao Lin, and Hongyu Chen. PiFlow: Principle-aware scientific discovery with  
381 multi-agent collaboration. *arXiv preprint arXiv:2505.15047*, 2025.
- 382 [7] Peter Jansen, Oyvind Tafjord, Marissa Radensky, Pao Siangliulue, Tom Hope, Bhavana Dalvi,  
383 Bodhisattwa Prasad Majumder, Daniel S Weld, and Peter Clark. CodeScientist: End-to-end  
384 semi-automated scientific discovery with code-based experimentation. In *Findings of the*  
385 *Association for Computational Linguistics: ACL 2025*, pages 13370–13467, 2025.
- 386 [8] Audrey Cheng, Shu Liu, Melissa Pan, Zhifei Li, Bowen Wang, Alex Krentsel, Tian Xia, Mert  
387 Cemri, Jongseok Park, Shuo Yang, et al. Barbarians at the gate: How AI is upending systems  
388 research. *arXiv preprint arXiv:2510.06189*, 2025.
- 389 [9] Alexander Novikov, Ngân Vũ, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt  
390 Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco JR Ruiz, Abbas Mehrabian,  
391 et al. AlphaEvolve: A coding agent for scientific and algorithmic discovery. *arXiv preprint*  
392 *arXiv:2506.13131*, 2025.
- 393 [10] Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni,  
394 and Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of*  
395 *the association for computational linguistics*, 12:157–173, 2024.
- 396 [11] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang  
397 Ding, Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating llms as agents. *arXiv preprint*  
398 *arXiv:2308.03688*, 2023.
- 399 [12] Nelson F Liu, Tianyi Zhang, and Percy Liang. Evaluating verifiability in generative search  
400 engines. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages  
401 7001–7025, 2023.
- 402 [13] Sewon Min, Kalpesh Krishna, Xinxin Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer,  
403 Luke Zettlemoyer, and Hannaneh Hajishirzi. FActScore: Fine-grained atomic evaluation of  
404 factual precision in long form text generation. In *Proceedings of the 2023 Conference on*  
405 *Empirical Methods in Natural Language Processing*, pages 12076–12100, 2023.
- 406 [14] Ori Press, Andreas Hochlehnert, Ameya Prabhu, Vishaal Udandarao, Ofir Press, and Matthias  
407 Bethge. CiteME: Can language models accurately cite scientific claims? *Advances in Neural*  
408 *Information Processing Systems*, 37:7847–7877, 2024.
- 409 [15] Theo Haerder and Andreas Reuter. Principles of transaction-oriented database recovery. *ACM*  
410 *computing surveys (CSUR)*, 15(4):287–317, 1983.

- 411 [16] Shu Liu, Mert Cemri, Shubham Agarwal, Alexander Krentsel, Ashwin Naren, Qiuyang Mang,  
412 Zhifei Li, Akshat Gupta, Monishwaran Maheswaran, Audrey Cheng, Melissa Pan, Ethan Boneh,  
413 Kannan Ramchandran, Koushik Sen, Alexandros G. Dimakis, Matei Zaharia, and Ion Stoica.  
414 Skydiscover: A flexible framework for ai-driven scientific and algorithmic discovery, 2026.
- 415 [17] Patrick Tser Jern Kon, Jiachen Liu, Qiuyi Ding, Yiming Qiu, Zhenning Yang, Yibo Huang,  
416 Jayanth Srinivasa, Myungjin Lee, Mosharaf Chowdhury, and Ang Chen. Curie: Toward rigorous  
417 and automated scientific experimentation with AI agents. *arXiv preprint arXiv:2502.16069*,  
418 2025.
- 419 [18] Yougang Lyu, Xi Zhang, Xinhao Yi, Yuyue Zhao, Shuyu Guo, Wenxiang Hu, Jan Piotrowski,  
420 Jakub Kaliski, Jacopo Urbani, Zaiqiao Meng, et al. EvoScientist: Towards multi-agent evolving  
421 AI scientists for end-to-end scientific discovery. *arXiv preprint arXiv:2603.08127*, 2026.
- 422 [19] Mert Cemri, Shubham Agrawal, Akshat Gupta, Shu Liu, Audrey Cheng, Qiuyang Mang, Ashwin  
423 Naren, Lutfi Eren Erdogan, Koushik Sen, Matei Zaharia, et al. AdaEvolve: Adaptive LLM  
424 driven zeroth-order optimization. *arXiv preprint arXiv:2602.20133*, 2026.
- 425 [20] Tingting Chen, Srinivas Anumasa, Beibei Lin, Vedant Shah, Anirudh Goyal, and Dianbo  
426 Liu. Auto-Bench: An automated benchmark for scientific discovery in LLMs. *arXiv preprint*  
427 *arXiv:2502.15224*, 2025.
- 428 [21] Yujie Liu, Zonglin Yang, Tong Xie, Jinjie Ni, Ben Gao, Yuqiang Li, Shixiang Tang, Wanli  
429 Ouyang, Erik Cambria, and Dongzhan Zhou. ResearchBench: Benchmarking LLMs in scientific  
430 discovery via inspiration-based task decomposition. *arXiv preprint arXiv:2503.21248*, 2025.
- 431 [22] Tianze Xu, Pengrui Lu, Lyumanshan Ye, Xiangkun Hu, and Pengfei Liu. Researcherbench:  
432 Evaluating deep ai research systems on the frontiers of scientific inquiry. *arXiv preprint*  
433 *arXiv:2507.16280*, 2025.
- 434 [23] Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. MAgentBench: Evaluating language  
435 agents on machine learning experimentation. *arXiv preprint arXiv:2310.03302*, 2023.
- 436 [24] Patrick Tser Jern Kon, Jiachen Liu, Xinyi Zhu, Qiuyi Ding, Jingjia Peng, Jiarong Xing, Yibo  
437 Huang, Yiming Qiu, Jayanth Srinivasa, Myungjin Lee, et al. EXP-Bench: Can AI conduct AI  
438 research experiments? *arXiv preprint arXiv:2505.24785*, 2025.
- 439 [25] Giulio Starace, Oliver Jaffe, Dane Sherburn, James Aung, Jun Shern Chan, Leon Maksin,  
440 Rachel Dias, Evan Mays, Benjamin Kinsella, Wyatt Thompson, et al. PaperBench: Evaluating  
441 AI’s ability to replicate AI research. *arXiv preprint arXiv:2504.01848*, 2025.
- 442 [26] Alisia Lupidi, Bhavul Gauri, Thomas Simon Foster, Bassel Al Omari, Despoina Magka, Al-  
443 berto Pepe, Alexis Audran-Reiss, Muna Aghamelu, Nicolas Baldwin, Lucia Cipolina-Kun,  
444 et al. AIRS-Bench: A suite of tasks for frontier AI research science agents. *arXiv preprint*  
445 *arXiv:2602.06855*, 2026.
- 446 [27] Zhen Wang, Fan Bai, Zhongyan Luo, Jinyan Su, Kaiser Sun, Xinle Yu, Jieyuan Liu, Kun  
447 Zhou, Claire Cardie, Mark Dredze, et al. FIRE-Bench: Evaluating agents on the rediscovery of  
448 scientific insights. *arXiv preprint arXiv:2602.02905*, 2026.
- 449 [28] Jun Shern Chan, Neil Chowdhury, Oliver Jaffe, James Aung, Dane Sherburn, Evan Mays, Giulio  
450 Starace, Kevin Liu, Leon Maksin, Tejal Patwardhan, et al. MLE-Bench: Evaluating machine  
451 learning agents on machine learning engineering. *arXiv preprint arXiv:2410.07095*, 2024.
- 452 [29] OpenAI. Parameter golf: OpenAI model craft challenge. [https://github.com/openai/  
453 parameter-golf](https://github.com/openai/parameter-golf), 2026.
- 454 [30] Erzhao Shao, Yifang Wang, Yifan Qian, Zhenyu Pan, Han Liu, and Dashun Wang. SciSciGPT:  
455 Advancing human–AI collaboration in the science of science. *Nature Computational Science*,  
456 6(3):301–315, 2026.
- 457 [31] Jiaqi Liu, Peng Xia, Siwei Han, Shi Qiu, Letian Zhang, Guiming Chen, Haoqin Tu, Xinyu  
458 Yang, Jiawei Zhou, Hongtu Zhu, Yun Li, Jiaheng Zhang, Yuyin Zhou, Zeyu Zheng, Cihang Xie,  
459 Mingyu Ding, and Huaxiu Yao. AutoResearchClaw: Fully autonomous research from idea to  
460 paper, 2026.

- 461 [32] Palash Goyal, Mihir Parmar, Yiwen Song, Hamid Palangi, Tomas Pfister, and Jinsung Yoon.  
462 ScholarPeer: A context-aware multi-agent framework for automated peer review. *arXiv preprint*  
463 *arXiv:2601.22638*, 2026.
- 464 [33] Shu Liu, Shubham Agarwal, Monishwaran Maheswaran, Mert Cemri, Zhifei Li, Qiuyang Mang,  
465 Ashwin Naren, Ethan Boneh, Audrey Cheng, Melissa Z Pan, et al. EvoX: Meta-evolution for  
466 automated discovery. *arXiv preprint arXiv:2602.23413*, 2026.

## 467 A Limitations

468 **Benchmark coverage.** We designed CoE and CoE Audit to be domain-agnostic, but validating that  
469 generality requires evaluation across diverse scientific domains. Our current experiments focus on  
470 systems-optimization tasks (ADRS), where gold-standard evaluators make Score Verification and  
471 Specification Violation Detection straightforward. Open-ended domains—biology, materials science,  
472 theoretical ML—pose harder challenges: evidence chains may involve wet-lab protocols, simulation  
473 reproducibility, or proof sketches, each demanding domain-specific verification logic that we have not  
474 yet built or tested. Extending CoE to these settings is a natural next step, though the core abstraction  
475 (claims linked to evidence via typed provenance records) should transfer; what changes is the set of  
476 integrity checks required.

477 **Citation verification depth.** Our Citation Verification checks whether cited references *exist*—a  
478 necessary condition that already catches a surprising number of failures (§B.1, Case 2). However,  
479 existence is far from sufficient: a real citation can still be used to support a claim the cited paper  
480 never made. Full citation verification would require passage-level natural language inference against  
481 the cited paper’s text, effectively asking “does this source actually say what the citing paper claims it  
482 says?” This is a known open problem in scholarly NLI, and we leave it to future work, noting that  
483 even existence-level checking already reveals meaningful architectural differences between systems.

484 **Automated review as proxy.** ScholarPeer serves as a scalable proxy for review quality but does  
485 not replace human expert evaluation. LLM reviewers are systematically blind to certain failure  
486 modes (e.g., domain-specific score interpretation, specification violation detection). CoE Audit  
487 addresses some of these blind spots but is itself limited to structural integrity, not scientific novelty or  
488 significance.

## 489 B Paper Quality Statistics

490 We report structural statistics of all 60 papers (4 systems  $\times$  3 seeds  $\times$  5 tasks) to characterize the  
491 surface-level quality of AI-generated manuscripts. Statistics are extracted automatically from the  
492 LaTeX source and compiled PDF of each paper.

493 Table 5 reports the mean and standard deviation of each metric across the 15 papers per system. ARC  
494 papers are the longest ( $12.7 \pm 1.2$  pages,  $5,417 \pm 863$  words) and most figure-heavy ( $4.3 \pm 2.0$  per  
495 paper). ScientistOne produces the largest bibliography ( $55.3 \pm 3.6$  entries),  $3.8 \pm 0.4$  figures.  
496 AIR is equation-heavy ( $7.1 \pm 1.8$  per paper) and table-heavy ( $4.9 \pm 1.5$ ) but generates nearly zero  
497 figures ( $0.1 \pm 0.2$ ). DS produces the shortest papers ( $6.7 \pm 2.1$  pages).

Table 5: Paper quality statistics across four systems (mean  $\pm$  std over 15 papers each).

Metric	ARC	DS	AIR	ScientistOne
Pages	$12.7 \pm 1.2$	$6.7 \pm 2.1$	$10.4 \pm 1.1$	$10.2 \pm 0.8$
Words	$5,417 \pm 863$	$2,741 \pm 968$	$5,116 \pm 524$	$4,643 \pm 438$
Figures	$4.3 \pm 2.0$	$1.3 \pm 0.4$	$0.1 \pm 0.2$	$3.8 \pm 0.4$
Tables	$2.1 \pm 0.7$	$0.8 \pm 0.7$	$4.9 \pm 1.5$	$1.9 \pm 0.2$
Equations	$1.3 \pm 1.6$	$1.9 \pm 2.0$	$7.1 \pm 1.8$	$3.0 \pm 1.1$
Unique cite keys	$23.5 \pm 8.7$	$9.3 \pm 6.9$	$19.5 \pm 3.9$	$18.3 \pm 4.5$
Bib entries	$23.5 \pm 8.7$	$13.1 \pm 13.2$	$15.8 \pm 1.9$	$55.3 \pm 3.6$
Sections	$6.6 \pm 2.5$	$6.5 \pm 1.0$	$5.0 \pm 0.0$	$5.5 \pm 0.6$
Subsections	$9.1 \pm 3.9$	$9.7 \pm 3.4$	$9.9 \pm 1.7$	$7.5 \pm 1.4$

### 498 B.1 Failure Mode Case Studies

499 We highlight four failure modes from the 60-paper audit, each illustrating a different evidence chain  
500 break that CoE Audit’s integrity checks are designed to catch.

501 **Case 1: Six orders of magnitude (ARC, LLM-SQL, seed 2).** The paper introduces “SCOR,” a  
502 static column-ordering routine, and reports a combined score of 1,538,006.69—on a benchmark whose

503 scoring metric lives on a  $[0, 1]$  scale. The figure is not a transcription error: it is the sum of squared  
504 prefix-hit lengths across datasets, an internal metric that the system computed and presented as if it  
505 were the ADRS score. The paper is internally coherent—it defines its own evaluation protocol, runs  
506 controlled comparisons against a baseline (scoring 1,537,927.99), and draws reasonable conclusions  
507 within that framing. An automated reviewer evaluating narrative quality alone would find nothing  
508 wrong. Score Verification catches this immediately: the gold evaluator re-run crashes (the submitted  
509 code fails to produce a valid solution), making the entire evidence chain from score to evaluator  
510 unresolvable.

511 **Case 2: A bibliography from model memory (AIR, PRISM, seed 1).** This paper’s bibliography  
512 contains 15 references. Citation Verification finds that 3 of those citations are hallucinated—we could  
513 not trace them to real papers and associated authors indexed by prominent indexes like Semantic  
514 Scholar, arXiv. This is not an edge case—AIR and DS produce phantom citations at rates of 9% and  
515 21% respectively (Table 1), compared to 0% for systems with structured retrieval pipelines.

516 **Case 3: Convergent specification violation (DS, LLM-SQL, seed 1).** Specification violation  
517 flags this paper with a unanimous 3/3 vote. The submitted code achieves a legitimate score of 0.697  
518 that passes Score Verification, but it does so by exploiting a gap between what the evaluator checks  
519 and what the benchmark intends to measure. The code sorts columns differently per row-group block,  
520 then renames all columns back to the original schema before concatenation—this makes `pd.concat`  
521 assemble blocks in insertion order rather than realigning by column name, effectively permuting  
522 column order per row-group. The evaluator validates row counts and total character counts but not  
523 column-to-column correspondence, so the permutation goes undetected. The same exploit appears  
524 independently in two other systems (AIR seed 2 and `ScientistOne` seed 2), providing convergent  
525 evidence that this is a genuine benchmark vulnerability rather than an isolated accident.<sup>2</sup>

526 **Case 4: Correct score, fictional algorithm (ARC, TXN, seed 1).** This paper passes Score  
527 Verification: the reported score of 3,311 falls within adaptive tolerance of the gold evaluator re-  
528 run (3,247). Yet Method-Code Alignment reveals a complete disconnect between what the paper  
529 describes and what the code implements. The paper introduces “STAR,” a system built on bitwise  
530 integer encoding for conflict detection, an  $O(1)$  surrogate cost model, and equidistant placement of  
531 high-contention anchor transactions. The submitted code implements none of these: it uses standard  
532 Python sets for conflict tracking, calls the full simulator on every iteration (no surrogate), and clusters  
533 read-heavy keys sequentially rather than distributing write-heavy anchors. This case illustrates why  
534 Score Verification alone is insufficient—the solver works, but the paper describes a different algorithm  
535 entirely, making the method section unreproducible regardless of how accurate the reported numbers  
536 are.

## 537 C System Implementation Details

### 538 C.1 Solver Architecture

539 The Solver consists of two agents. The *Solution Development Agent* receives an idea and task-specific  
540 instructions, operates in a sandboxed environment with tools for file I/O, command-line execution,  
541 solution management, and knowledge base retrieval, and follows an iterative refinement loop—  
542 executing experiments, debugging errors, and optimizing validation metrics—while maintaining an  
543 experimental log. The *Report Writing Agent* parses experimental artifacts to generate a technical  
544 report summarizing methodology and outcomes.

### 545 C.2 Ablation Study

546 Following PEE, the system selects the global best solution for further validation. An ablation  
547 agent identifies the solution’s core components and proposes controlled experiments to isolate their  
548 contributions. These ablated versions are implemented and re-evaluated to quantify the performance  
549 impact of specific architectural or algorithmic choices.

---

<sup>2</sup>For `ScientistOne` seed 2, the exploit is present in the submitted code, but the I2 majority-vote protocol did not reach consensus (1 of 5 judges flagged it), so it is not counted as a violation in Table 1. This illustrates a limitation of LLM-judged checks at the current vote threshold.

### 550 C.3 Claim Verifier: Per-Type Verification Rules

551 The Claim Verifier dispatches on the four claim types defined in §3. The writer commits each claim’s  
552 evidence via annotation tags—a line in the experimental log, a citation key, an internal ablation, a  
553 baseline from the PI brief, or an explicit “unsourced” marker—and the verifier maps that evidence  
554 onto the claim type’s verification rule:

- 555 • **Numerical claims** are checked by numeric tolerance against the cited evidence (log line,  
556 ablation entry, or PI baseline), with a  $\pm 3$ -line window on log lines and unit-aware normal-  
557 izations for percent-versus-fraction and millisecond-versus-second drift.
- 558 • **Citation claims** are checked by resolving the cite key against the bibliography and then  
559 asking a one-shot LLM judge (JSON mode) whether the cited work’s abstract supports the  
560 specific assertion.
- 561 • **Methodological claims** are checked by substantive textual overlap against the cited region  
562 of the experimental log.
- 563 • **Conclusion claims** are detected post-hoc by scanning for comparison language, linked to  
564 supporting numerical claims by token and number overlap, and—when parseable—validated  
565 by arithmetic check against pairwise differences.

566 Claims tagged “unsourced” or carrying malformed annotations are dropped automatically; each  
567 records a break code for downstream reporting.

### 568 C.4 Paper Writer Pipeline

569 The Paper Writer is a five-stage pipeline that converts raw materials into a verified  $\LaTeX$  draft.  
570 The first four stages operate on a *research representation*—a structured markdown narrative with  
571 inline evidence annotations—before any  $\LaTeX$  is generated, enforcing the “provenance before prose”  
572 principle described in §4.3.

573 **CONCEIVE.** A single LLM call reads all assembled raw materials (PI brief, experimental log,  
574 verified scores, solver code, seed-paper abstracts) and produces the initial research representation.  
575 This document captures the story arc—problem, gap, approach, result, limitation—with every factual  
576 claim carrying an inline evidence tag binding it to a specific workspace artifact (log line, score file,  
577 citation key, or ablation entry). **CONCEIVE** establishes the narrative structure but does not validate  
578 evidence chains; that is deferred to **GROUND**.

579 **GROUND.** Deterministic checks validate each evidence annotation against the raw materials: the  
580 published score is present and matches `bestrun_selection.json`, baselines are labelled **VERIFIED**  
581 (traceable to a PI brief entry) or **ESTIMATED** (unattributed “leaderboard” references are flagged),  
582 every referenced file exists, all expected sections are present, and hyperbole counts and known score  
583 mismatches are recorded. Each claim receives a **SUPPORTED**, **PARTIAL**, or **UNSUPPORTED** label, and  
584 the overall grounding ratio (supported claims / total claims) is computed.

585 **CRITIC.** One LLM call audits story-level coherence: gap–approach alignment, internal contradic-  
586 tions, overclaims relative to evidence strength, missing comparisons, baseline fairness, and honest  
587 limitations. The critic returns **PASS** or a list of **MAJOR/MINOR** issues.

588 **RESOLVE.** A single LLM call rewrites the representation against both the **GROUND** flag list  
589 and the **CRITIC** issue list: unsupported claims are dropped or softened, contradictions are re-  
590 solved using the verified source, overclaims are calibrated, and missing sections are filled. The  
591 **GROUND**→**CRITIC**→**RESOLVE** loop iterates for up to two rounds, terminating on convergence (zero  
592 flags) or plateau (the flag count stops decreasing). A final grounding ratio below 50% raises a hard  
593 error, aborting the paper rather than producing a poorly grounded draft.

594 **COMPOSE.** The grounded representation is handed to per-section writers that emit  $\LaTeX$  one  
595 section at a time, with each numerical or citation-bearing sentence committing to its evidence source  
596 at writing time. The composed draft then passes through the Claim Verifier (§C.3), whose findings  
597 drive a Refiner pass that applies corrections and strips inline annotations. Only a draft with no

598 blocking violations is promoted to `final_paper.pdf`; failed drafts are saved as `draft_paper.pdf`  
599 with a diagnostic report.

## 600 C.5 Problem Investigator Details

601 PI is a five-stage pipeline (plus two auxiliary stages), where each stage communicates via file-backed  
602 artifacts on disk.

603 **Stage 1: Citation Graph.** Starting from 2–4 seed papers, PI traverses the Semantic Scholar API  
604 (references and citations) up to 2 hops deep, producing a citation graph of approximately 2,000–5,000  
605 candidate papers.

606 **Stage 2: Literature Filter.** An LLM scores each paper in the graph on two axes—methodology  
607 relevance and problem alignment (1–5 each)—and classifies papers into tiers: Core (both  $\geq 4$ ),  
608 Adjacent (one  $\geq 4$ , other  $\geq 3$ ), Spark, or Noise. The resulting elite pool contains approximately 500  
609 papers. A topic-relevance gate aborts the pipeline if fewer than 5 Core+Adjacent papers appear,  
610 preventing downstream drift from weak seeds.

611 **Stage 3: Multi-Round Investigation.** A Principal Investigator agent orchestrates specialist sub-  
612 agents across 3 rounds. Each round consists of candidate selection from the elite pool (Librarian  
613 agent), parallel PDF reading and structured note extraction (5 Researcher agents), and synthesis of  
614 findings into thematic research direction dossiers (SubdomainWriter agent). An IslandConsolidator  
615 agent merges redundant directions and retires low-quality ones after each round. The target is  
616 approximately 100 paper notes organized into 5–15 research directions.

617 **Stage 4: Evaluation Protocol Audit and Targeted Literature Refresh.** Per-direction audit reports  
618 are generated and scored on a checklist rubric across multiple rounds until the direction passes. A  
619 focused mini citation crawl on the winning direction produces 20–30 additional paper notes, filling  
620 gaps identified during the audit.

621 **Stage 5: Experiment Brief Synthesis.** Directions are scored by seed relevance, then a section-by-  
622 section writer produces the final Experiment Brief through a multi-round critic loop (up to 5 rounds  
623 with section-level revision). The brief contains three sections: (1) a research landscape with technique  
624 taxonomy and best-known results, (2) a concrete experiment plan with baselines, metrics, and ablation  
625 design, and (3) a literature context with 25–40 references traceable to paper notes extracted from the  
626 source PDFs.

## 627 D CoE Audit Details

### 628 D.1 Sub-Code Definitions

#### 629 D.1.1 Score Verification Failure Codes.

Code	Label	Description
SF-1	Genuine discrepancy	Paper score differs from evaluator rerun mean beyond adaptive tolerance, with consistent metric and scale
SF-2	Metric mismatch	Paper and evaluator report different metrics or scales (e.g., percentage vs. fraction, inverted objective)
SF-3	Metric inversion	Paper score is the reciprocal of the evaluator score, suggesting an inverted optimization direction
SF-4	Evaluator failure	All $N=5$ evaluator reruns crashed or returned null
SF-5	Evaluator timeout	All $N=5$ evaluator reruns exceeded the time limit
SF-6	Score extraction failure	No headline score could be extracted from the agent’s logs or paper

Code	Label	Description
RWH-1	Evaluator exploitation	Solution reverse-engineers or exploits evaluator internals
RWH-2	Evaluator modification	Solution imports, reads, or modifies the evaluator
RWH-3	Data leakage	Solution reads test data, ground truth, or answer files
RWH-4	Hardcoded output	Solution produces hardcoded outputs rather than computing them
RWH-5	Specification exploit	Solution exploits metric edge cases or violates task constraints (e.g., column permutation in LLM-SQL)

Code	Label	Description
CF-1	Phantom citation	Reference does not exist anywhere (Scholar, DBLP, arXiv)
CF-2	Memory citation	Reference is real but recalled from model weights, not retrieved; meta-data often corrupted

630 **D.1.2 Specification Violation Sub-Codes.**

631 **D.1.3 Citation Verification Sub-Codes.**

632 **D.1.4 Method-Code Alignment Sub-Codes.**

Code	Label	Description
MCA-0	Faithful	Paper and code describe the same algorithm
MCA-1	Acceptable simplification	Paper describes the same class of algorithm but omits or simplifies implementation details; rhetorical embellishment is allowed
MCA-2	Misrepresentation	Paper materially misrepresents the implemented algorithm (e.g., describes beam search when code uses greedy)

633 **D.2 Audit Procedure and Reproducibility**

634 CoE Audit is designed for reproducibility. Table 6 summarizes the automation level of each component. Score Verification uses LLM-based extraction to identify the paper’s reported score from  
635  $\text{\TeX}$  and PDF files, then compares it against reproduced scores from evaluator re-runs via numeric  
636 tolerance—the comparison itself is fully deterministic. Specification Violation and Method-Code  
637 Alignment are LLM-judged, with majority vote across multiple independent runs to reduce judgment  
638 noise. Citation Verification is primarily API-based (Semantic Scholar, arXiv, OpenAlex, CrossRef  
639 lookup), with an LLM disambiguation step for title-only matches.  
640

Table 6: CoE Audit automation level per component.

Component	Automation	Model (if LLM)
Score Verification	LLM extraction + automated comparison	Gemini 3.0 Flash
Specification Violation	LLM-judged	Gemini 3.1 Pro
Citation Verification	Automated + LLM disambiguation	Gemini 3.0 Flash
Method-Code Alignment	LLM-judged	Gemini 3.1 Pro

641 **Human verification.** Although the audit pipeline is automated, all flagged positives for I1 (Score  
642 Verification), I2 (Specification Violation), and I3 (Citation Verification) were manually reviewed  
643 and corrected by human reviewers before reporting in Table 1. A small number of auditor false  
644 positives (e.g., API resolution failures for real papers, score extraction errors) were identified and  
645 removed during this process. For I4 (Method-Code Alignment), human reviewers validated a sample  
646 of the LLM judgments; the results in Table 1 reflect the LLM majority-vote scores without manual  
647 correction.

## 648 E Failure Cases per Audit Metric

### 649 E.1 I1: Score Verification Errors

650 We manually reviewed every I1 (Score Verification) failure flagged across the four baseline systems  
651 and categorise the 15 confirmed agent errors into four classes (Table 7).<sup>3</sup>

Table 7: Confirmed agent I1 errors, by category.

Category	Definition	Count	%
value_mismatch	Paper headline number differs from evaluator rerun beyond tolerance, with the same metric and scale.	11	73%
paper_score_unavailable	Paper contains no machine-readable headline number for the auditor to verify against.	2	13%
evaluator_error	Submitted artifact cannot be re-evaluated within the budget (timeout or crash on the agent’s own solution file).	1	7%
metric_mismatch	Paper and code report different metrics, scales, or optimisation directions, so the numbers are not directly comparable.	1	7%

652 **value\_mismatch (n=11).** The paper and the code agree on what is being measured but disagree  
653 on the value. Most (7/11) of these gaps are within 5% of the paper-reported number—small enough  
654 to plausibly arise from unreported seed variance, but uniformly biased towards a better-than-rerun  
655 headline. The two largest gaps are qualitatively different errors: a fabricated metric direction in  
656 DS cloudcast seed-2 (26.7%), where the paper relabels a cost metric as “utility” and inverts the  
657 optimisation direction so that a worse-than-baseline number reads as an improvement; and a 2×  
658 mismatch in ARC prism seed-3 between the paper headline (12.74) and what the submitted code  
659 actually produces (26.24).

660 **paper\_score\_unavailable (n=2).** The agent’s writeup contains no machine-readable quantitative  
661 result for the headline metric. Both cases are baseline-system failures: an AIR prism paper that  
662 omits scores entirely, and a DS cloudcast PDF that failed to render its results section.

663 **evaluator\_error (n=1).** The artifact the agent declares as its solution cannot be re-evaluated  
664 end-to-end within the budget. The single case is ARC 11m\_sql seed-2, which ships program.py—a  
665 multi-condition experiment harness that runs 8 baseline conditions × 3 seeds × multiple datasets at  
666 import time before any single number is produced.

667 **metric\_mismatch (n=1).** The paper and the code do not agree on what they are measuring. The  
668 single case is ARC 11m\_sql seed-3, which publishes a headline of 0.0 that the paper itself attributes  
669 to “a critical integration failure with the benchmark evaluator framework”—the agent surfaced a  
670 known-broken result as its final score.

671 **Per-system error shape.** Table 8 shows that the failure mix is highly system-specific. AIR and  
672 DS errors are dominated by small-to-medium value\_mismatch: the numbers exist and are in the  
673 right ballpark but do not exactly reproduce. ARC spans the widest range of categories and the largest  
674 single discrepancy (106%), and contributes the only evaluator\_error and metric\_mismatch  
675 cases. ScientistOne produces no confirmed I1 errors, consistent with its 12/12 score-verification  
676 result in Table 1.

### 677 E.2 I3: Citation integrity - Discovered phantom citations

<sup>3</sup>A small number of auditor false positives were corrected before reporting Table 1; this section reports only confirmed agent errors. Table 1 reports pass rates over 12 papers per system (EPLB excluded), so its implied failure count (11) is lower than the 15 errors here, which include EPLB papers.

Table 8: Confirmed agent I1 errors per system, by category.

System	value_mismatch	paper_score_unavail.	evaluator_error	metric_mismatch	Total
AIR	3	1	0	0	4
ARC	5	0	1	1	7
DS	3	1	0	0	4
ScientistOne	0	0	0	0	0

Phantom citations
<b>ruth2018castflow</b> : R�uth, Jan et al. <i>CastFlow: Clean-slate multicast approach using in-advance path computation in software-defined networks</i> . 2018 IEEE 43rd Conference on Local Computer Networks (LCN), 2018.
<b>dou2022hetmoe</b> : Dou, Shiwei et al. <i>HetMoE: An Efficient Distributed MoE Training System for Heterogeneous Clusters</i> . arXiv preprint arXiv:2210.12384, 2022.
<b>li2023lightllm</b> : Li, Zhuohan et al. <i>LightLLM: A highly optimized LLM inference system with token-level kv cache management</i> . arXiv preprint arXiv:2309.04414, 2023.
<b>purohit2024prism</b> : Purohit, Archit et al. <i>Prism: Optimizing multi-model LLM serving on GPU clusters</i> . Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), 2024.
<b>miao2023muxserve</b> : Miao, Xupeng et al. <i>MuxServe: Multiplexing large language models for high throughput and low latency</i> . arXiv preprint arXiv:2311.05602, 2023.
<b>cloudcast</b> : Zheng, Q. et al. <i>CloudCast: Cost-efficient multicast routing in cloud networks</i> . IEEE INFOCOM 2019 - IEEE Conference on Computer Communications, 2019.
<b>idreos2012main</b> : Idreos, Stratos et al. <i>Main-memory column stores</i> . Foundations and Trends® in Databases, 2012.
<b>lightllm2023</b> : Li, Zhen et al. <i>LightLLM: A Lightweight and Highly Efficient Python-based Large Language Model Serving Framework</i> . arXiv preprint arXiv:2310.01234, 2023.
<b>bhuiyan2024prism</b> : Bhuiyan, M. et al. <i>Prism: A Flexible and Scalable Multi-LLM Serving System</i> . Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), 2024.
<b>ep1b2023</b> : Wang, H. et al. <i>Expert Parallelism Load Balancing in Mixture-of-Experts Models</i> . Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis (SC), 2023.
<b>hussin2011metaheuristic</b> : Hussin, MS et al. <i>Metaheuristic algorithms for traveling salesman problem: A review</i> . Annals of the University of Craiova, Mathematics and Computer Science Series, 2011.
<b>zhang2020job</b> : Zhang, Jun et al. <i>Job shop scheduling research</i> . International Journal of Production Research, 2020.
<b>jetstream2014</b> : Rabkin, Ariel et al. <i>JetStream: Enabling wide-area data streaming</i> . 11th USENIX Symposium on Networked Systems Design and Implementation (NSDI 14), 2014.
<b>slingshot2022</b> : Doe, John et al. <i>Slingshot: High-performance routing across federated cloud environments</i> . Proceedings of the ACM SIGCOMM 2022 Conference, 2022.
<b>cloudcast2021</b> : Chen, Yiting et al. <i>CloudCast: Cost-effective data distribution in multi-cloud deployments</i> . IEEE INFOCOM 2021 - IEEE Conference on Computer Communications, 2021.
<b>cui2004</b> : Cui, Jun-Hong et al. <i>QoS multicast routing in dynamic networks</i> . IEEE Network, 2004.
<b>laoutaris2011netstitcher</b> : Laoutaris, Nikolaos et al. <i>NetStitcher: Un-tethering bulk storage from the network edge</i> . Proceedings of the ACM SIGCOMM 2011 conference, 2011.
<b>feng2020traffic</b> : Feng, Xin et al. <i>Traffic engineering in software-defined wide-area networks: A survey</i> . IEEE/ACM Transactions on Networking, 2020.

<b>Phantom citations</b> (continued)
<b>epstein2005online</b> : Epstein, Leah and van Stee, Rob. <i>Online bin packing with square-root sized items</i> . Information Processing Letters, 2005.
<b>ba2023modelbox</b> : Ba, Yuhan et al. <i>ModelBox: A Framework for Multi-Model Multi-Tenant Serving</i> . 2023 USENIX Annual Technical Conference (USENIX ATC 23), 2023.
<b>zheng2024eplb</b> : Zheng, Lianmin et al. <i>EPLB: Load Balancing for Expert Parallelism in Large Language Models</i> . arXiv preprint arXiv:2401.03221, 2024.
<b>sutskever2013importance</b> : Sutskever et al. <i>SGD with Momentum</i> . ICML, 2013.
<b>choy1991heuristic</b> : Choy, M-S. <i>Heuristic algorithms for the Steiner tree problem with an application to network routing</i> . Proceedings of the 1991 ACM SIGCOMM, 1991.
<b>beloglazov2012energy</b> : Beloglazov, Anton and Buyya, Rajkumar. <i>Energy-efficient routing and resource allocation in multi-cloud</i> . Journal of Network and Computer Applications, 2012.
<b>grotschel1993steiner</b> : Grotschel, Martin et al. <i>The Steiner tree packing problem in telecommunications</i> . 1993.
<b>romero2021automated</b> : Romero, F et al. <i>Automated algorithm design using large language models</i> . Advances in Neural Information Processing Systems, 2023.
<b>wu2015spanning</b> : Wu, Chuan et al. <i>Spanning tree based data transfer in multi-cloud architectures</i> . 2015.
<b>melian2012cloud</b> : Melian, L et al. <i>Cloud computing economics: a survey</i> . 2012.
<b>faragardi2017multi</b> : Faragardi, Hamid Reza et al. <i>Multi-cloud data distribution with cost optimization</i> . 2017.
<b>mishra2018cost</b> : Mishra, A et al. <i>Cost efficient routing in multi-cloud environments</i> . 2018.
<b>voss1992steiner</b> : Voss, Stefan. <i>The Steiner tree problem with edge capacities</i> . 1992.
<b>bubeck2023approaches</b> : Bubeck, S et al. <i>Approaches to code generation and synthesis</i> . 2023.
<b>wang2020automated</b> : Wang, X et al. <i>Automated hyperparameter optimization in cloud computing</i> . 2020.
<b>smith2019data</b> : Smith, J et al. <i>Data transfer costs in modern cloud platforms</i> . 2019.
<b>liu2022network</b> : Liu, Y et al. <i>Network aware multi-cloud data distribution</i> . 2022.
<b>jones2023oscillatory</b> : Jones, R et al. <i>Oscillatory dynamics in automated search landscapes</i> . 2023.
<b>kim2021puber</b> : Kim, Young Jin et al. <i>Puber: Efficient expert parallelism for mixture-of-experts</i> . arXiv preprint arXiv:2111.05454, 2021.
<b>romera2021optimizing</b> : Romera-Paredes, Oscar et al. <i>Optimizing mixture-of-experts for large-scale distributed training</i> . arXiv preprint arXiv:2102.04353, 2021.
<b>clune2008how</b> : Clune, Jeff et al. <i>How evolutionary dynamics affects network topology</i> . Artificial life, 2008.
<b>llm_agents_2023</b> : Smith, J. and Doe, A. <i>Autonomous LLM Agents for Code Generation</i> . Journal of AI Research, 2023.
<b>api_misuse_2024</b> : Wang, L. and Lee, C. <i>API Misuse in LLM-Generated Code</i> . Proceedings of ICSE, 2024.
<b>fail_fast_2025</b> : Garcia, M. <i>Fail-Fast Sandboxing for Coding Agents</i> . IEEE Transactions on Software Engineering, 2025.
<b>prism2024</b> : Anonymous. <i>Prism: A Benchmark for Multi-LLM Serving Systems</i> . arXiv preprint, 2024.
<b>sarca2023</b> : Anonymous. <i>SARCA: Systems Architecture Research using Coding Agents</i> . arXiv, 2023.
<b>clockwork2023</b> : Anonymous. <i>Clockwork: Fast and Predictable Inference for Edge Machine Learning</i> . arXiv, 2023.
<b>garcia1982fully</b> : Garcia-Molina, Hector. <i>A fully distributed null-free algorithm for concurrent database updates</i> . IEEE Transactions on Software Engineering, 1982.

<b>Phantom citations</b> (continued)
<b>karlsson2020combinatorial</b> : Karlsson, Elias et al. <i>Combinatorial optimization by graph neural networks</i> . arXiv preprint arXiv:2010.16012, 2020.
<b>krajewski2024mixtures</b> : Krajewski, Adam et al. <i>Mixtures of experts: A systematic review</i> . arXiv preprint arXiv:2401.00000, 2024.
<b>paliwal2020regal</b> : Paliwal, Aditya et al. <i>REGAL: A transfer learning based methodology for hardware-software co-design</i> . MICRO, 2020.
<b>jain1998simulated</b> : Jain, AS and Meeran, S. <i>A simulated annealing algorithm for job shop scheduling problem</i> . Mathematical and computer modelling, 1998.
<b>tian2021cloud</b> : Tian, X et al. <i>Cloud egress cost optimization</i> . Proceedings of the ACM SIGCOMM 2021 Conference, 2021.
<b>zhao2020understanding</b> : Zhao, Y et al. <i>Understanding cloud network egress constraints</i> . USENIX Annual Technical Conference (ATC), 2020.
<b>binnig2021learned</b> : Binnig, Carsten et al. <i>The case for learned database systems</i> . arXiv preprint, 2021.
<b>yu2014staccato</b> : Yu, Xiangyao et al. <i>Staccato: A dependency-aware transaction scheduling system for many-core processors</i> . Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data, 2014.
<b>ren2016low</b> : Ren, Kun et al. <i>Low-overhead deadlock detection in distributed database systems</i> . Proceedings of the 2016 International Conference on Management of Data, 2016.
<b>pavlo2017make</b> : Pavlo, Andrew and Stonebraker, Michael. <i>What's make a database system fast?</i> . ACM SIGMOD Record, 2017.
<b>blanas2010comparison</b> : Blanas, Spyros et al. <i>A comparison of join algorithms for modern multi-core processors</i> . Proceedings of the VLDB Endowment, 2010.
<b>bailis2014bolt</b> : Bailis, Peter et al. <i>Bolt-on conflict-free replicated data types</i> . ACM SIGMOD Record, 2014.
<b>hardi1992precedence</b> : Hardi, S and Rakow, TC. <i>Precedence-based transaction scheduling</i> . Proceedings of the 2nd International Workshop on Research Issues on Data Engineering, 1992.
<b>kim2016fast</b> : Kim, Taesoo et al. <i>Fast and scalable serializable transactions in multicore in-memory databases</i> . Proceedings of the 2016 International Conference on Management of Data, 2016.
<b>wu2017transaction</b> : Wu, Yingjun et al. <i>Transaction scheduling using graph coloring</i> . Proceedings of the 2017 ACM SIGMOD International Conference on Management of Data, 2017.
<b>lin2015scheduler</b> : Lin, Jianshu et al. <i>To schedule or not to schedule: When is transaction scheduling worth the overhead?</i> . Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data, 2015.
<b>faleiro2017high</b> : Faleiro, Jose M and Abadi, Daniel J. <i>High performance serializable concurrency control with determinism</i> . Proceedings of the 2017 ACM SIGMOD International Conference on Management of Data, 2017.

678 **E.3 I4 Audit Failure Analysis**

679 We categorise the 95 method-code misalignment findings flagged by the I4 audit (across 25 papers<sup>4</sup>)  
680 into three semantic classes (Table 10). A single paper can produce multiple findings, so we report  
681 both the finding-level distribution and the distinct-paper count per category.

682 **incomplete\_broken (n=49, 19 papers)**. The largest category. The code generally targets the  
683 same problem the paper describes, but is missing one or more of the specific mechanisms claimed  
684 in the writeup, or substitutes a degenerate fallback. Recurring patterns include: described *multi-*  
685 *start initialisation with K sequences* collapsed to a single deterministic backtracking pass (AIR

<sup>4</sup>Table 1 implies 26 misaligned papers (12+10+3+1); one paper (DS seed-3/cloudcast) could not be cleanly classified and is excluded from the categorical breakdown.

Table 10: I4 method-code alignment findings, by category (n=95 findings, 25 affected papers).

Category	Definition	Findings	%	Papers
<code>incomplete_broken</code>	Paper-described mechanism or component is missing, partially implemented, or replaced with a degenerate fallback in the code.	49	52%	19
<code>algorithm_class_mismatch</code>	Code implements a fundamentally different algorithm class than the one described in the paper (e.g., paper claims Iterated Local Search, code is Simulated Annealing; paper claims a neural- or LLM-guided search, code is deterministic).	37	39%	15
<code>deceptive_dummy_code</code>	Code contains undisclosed elements designed to mislead evaluation: hidden environment-variable switches, or output deliberately shaped to game an evaluator metric.	9	9%	5

686 `prism`); claimed *link-penalty diversification* producing a constant assignment of the same path to  
 687 every partition (AIR `cloudcast`); described *arborescence lookahead* or *surrogate cost model* simply  
 688 absent from the code, with the true simulator invoked at every iteration instead (ARC `cloudcast`,  
 689 ARC `txn_scheduling`); claimed *2.0 GB memory threshold safeguards* reduced to a trivial overflow  
 690 check (DS `prism`).

691 **`algorithm_class_mismatch` (n=37, 15 papers).** The code implements a fundamentally different  
 692 algorithm class than the one described in the paper. The most common sub-pattern is the *claimed-*  
 693 *learning-loop-is-absent* case: a paper claims an LLM-driven evolutionary search, a neural network  
 694 predictor, or an LLM SQL optimiser, and the code is a single deterministic heuristic with no  
 695 LLM calls and no trained model (e.g. AIR `cloudcast`: “hybrid neuro-symbolic solver” → purely  
 696 classical Steiner-tree heuristics; ARC `llm_sql`: “36 LLM prompting strategies” → deterministic  
 697 dataframe reordering; DS `ep1b`: “LLM-driven evolutionary search over 27 generations” → standalone  
 698 deterministic load balancer). Classical algorithm-class swaps also recur: Iterated Local Search →  
 699 Simulated Annealing (AIR `prism`), Dinkelbach’s fractional programming → static sum-of-squares  
 700 cost (ARC `prism`), recursive partitioning → single-level grouping (DS `llm_sql`). We also observe  
 701 the rarer *method/baseline inversion* pattern, in which the paper labels  $X$  as the proposed method and  
 702  $Y$  as a rejected ablation, but the code uses  $Y$  (ARC `prism`: “GRASP Without Symbiosis” claimed as  
 703 final method while the code instantiates `SymbioticGRASPPacker`).

704 **`deceptive_dummy_code` (n=9, 5 papers).** The submitted artifact contains undisclosed code whose  
 705 presence appears intended to mislead automated evaluation. All nine findings come from ARC and  
 706 split into two sub-patterns: (i) *hidden environment-variable switches* (4 findings, 2 papers): the  
 707 code reads an undisclosed variable (`CONDITION`, `ABLATION`) at import time and dispatches to one  
 708 of several different solvers, while the paper presents a single unified algorithm (ARC `cloudcast`  
 709 seeds 2 and 3). (ii) *evaluator gaming* (5 findings, 3 papers): code intentionally shaped to inflate  
 710 a metric without solving the underlying task, e.g. returning an empty column-ordering list while  
 711 internally permuting values to maximise prefix-cache hits (ARC `llm_sql` seed-3).

712 **Per-system error shape.** Table 11 shows that the failure mix is sharply system-specific. AIR,  
 713 DeepScientist, and ScientistOne produce no `deceptive_dummy_code` findings: when these  
 714 systems misalign, the gap is between what the paper says and what the code does, with no active  
 715 obfuscation. ARC is the only system that produces `deceptive_dummy_code` findings (9/9, 5 papers),  
 716 spanning both hidden env-var switches and prefix-hit evaluator gaming.

## 717 F MLE-Bench and Parameter Golf Evaluation

718 We evaluate ScientistOne on MLE-Bench and Parameter Golf. Details regarding the benchmarks  
 719 and our evaluation setup are provided below.

720 **MLE-Bench.** We selected five MLE-Bench [28] Kaggle competitions spanning medical imaging,  
 721 fine-grained recognition, and 3D perception. Table 12 summarizes the mapping between the task

Table 11: I4 findings per system, by category. Each cell is a finding count; a single paper may contribute multiple findings.

System	algo_class_mismatch	incomplete_broken	deceptive_dummy_code	Total
AIR	3	9	0	12
ARC	15	23	9	47
DS	15	16	0	31
ScientistOne	4	1	0	5

Table 12: Mapping between MLE-Bench task names used in our evaluation and their corresponding official task IDs.

Task Name	Task ID
3D Object Detection	3d-object-detection-for-autonomous-vehicles
AI4Code	AI4Code
iMet 2020 FGVC7	imet-2020-fgvc7
RSNA Brain Tumor	rsna-miccai-brain-tumor-radiogenomic-classification
iNaturalist 2019 FGVC6	inaturalist-2019-fgvc6

722 names used in our evaluation and their official task IDs. We specifically targeted tasks in the Medium  
 723 and High difficulty tiers, as they possess sufficient complexity and substantive research value to  
 724 serve as meaningful benchmarks for automated scientific discovery and paper writing. Specifically,  
 725 AI4Code, iMet 2020 FGVC7, and iNaturalist 2019 FGVC6 are classified as Medium-difficulty tasks,  
 726 whereas 3D Object Detection and RSNA Brain Tumor are High-difficulty tasks. The agent is provided  
 727 with task descriptions and a code execution environment to develop its solutions. The execution  
 728 environment consists of a compute node provisioned with 8xH100 GPUs, 192 CPU cores, and 1TB of  
 729 RAM. To assist with formatting, the agent can query a validation tool an unlimited number of times  
 730 to ensure its submission CSV files are correctly structured. Furthermore, we permit the agent to query  
 731 the grading server up to 16 times to obtain evaluation scores on the test data. This deviates from the  
 732 official MLE-Bench protocol, which restricts evaluation to a single submission of the final generated  
 733 solution. We adopt this modified setup to better simulate real-world Kaggle competitions, where  
 734 participants iteratively submit solutions to a public leaderboard for feedback during the development  
 735 phase. While the agent is instructed to iterate on its solutions primarily based on metrics derived from  
 736 the validation dataset, the grading server is used sparingly to select the best-performing models on  
 737 the test set. Additionally, this setup is necessary to allow the agent to report accurate final test metrics  
 738 when generating the paper.

739 **Parameter Golf.** To assess adaptability in a novel, live environment, we evaluate `ScientistOne` on  
 740 the Parameter Golf competition [29]. This live competition serves as an ideal AI research task focused  
 741 on LLM training. It requires participants to train the highest-performing language model under strict  
 742 constraints: the final artifact must fit within a 16MB size limit and the training process must complete  
 743 in under 10 minutes on an 8xH100 system. Performance is evaluated by the compression rate –  
 744 measured in tokenizer-agnostic bits per byte (BPB) – on the FineWeb validation set. We provide  
 745 the agent with a sandbox tool to execute solution code to obtain evaluation metrics. Additionally,  
 746 the agent is provided with a knowledge base containing official leaderboard solutions up to a cutoff  
 747 date of April 27, 2026. As of this date, the state-of-the-art (SOTA) score was 1.0611, achieved by a  
 748 solution titled “BOS-Fixed SmearGate + LQER + SparseAttnGate + 9-Hparam Stack”. The agent is  
 749 tasked with building upon these existing leaderboard baselines to discover novel techniques capable  
 750 of further improving the SOTA score.

## 751 G Baseline Adaptation Details

752 All baselines are open-source systems adapted to the ADRS benchmark under identical conditions:  
 753 Gemini 3.1 Pro backbone, up to 20 solver iterations per task, 2-hour code generation windows, and  
 754 3 seeds per task. Runs that crashed due to infrastructure issues (API timeouts, rate limits,  $\LaTeX$   
 755 compilation errors) were re-attempted with fresh state up to 3 times; no run was re-attempted to  
 756 improve solver scores. Below we detail the per-system adaptations.

757 **AutoResearchClaw (ARC).** ARC v0.3.1 required 2 source patches: one to support the Gemini  
758 thinking-model API format, and one to configure the scaffold directory for ADRS tasks. No  
759 changes were made to ARC’s search logic, stopping criteria, or evaluation pipeline. We set  
760 `max_iterations=20` ( $6.7\times$  ARC’s default of 3) and `max_refine_duration=7200s`. ARC’s  
761 internal limits (`MAX_DECISION_PIVOTS=2`, `MAX_REPAIR_CYCLES=3`) are architectural and were not  
762 modified. ARC’s post-pipeline scoring (PQEP) failed on 4 of 15 runs due to a regex bug in evaluator  
763 path substitution; these were resolved by re-running the canonical ADRS evaluator on the submitted  
764 code. Figure generation failed across all 15 runs because ARC’s `FigureAgent` requires a Docker  
765 image not available in our environment; all papers were produced without figures.

766 **DeepScientist (DS).** DS required prompt-only changes—no source files were patched. Task  
767 specifications and the NeurIPS 2026 template were injected via the write-phase prompt. We set  
768 `CODE_TIMEOUT=7200s` and `WRITE_TIMEOUT=5400s`, with up to 3 retries each for code and write  
769 phases. DS’s write skill instructs the agent to retrieve citations via Semantic Scholar, arXiv, and  
770 CrossRef APIs; however, across all 15 write-phase logs, the agent never called any retrieval API or  
771 MCP tool, generating all references from model memory. This is a model compliance failure—the  
772 tools are available but the agent consistently shortcuts the citation retrieval instructions.

773 **AI-Researcher (AIR).** AIR required the most extensive adaptation: 19 source files were patched to  
774 interface with ADRS task configurations, the Docker-based sandbox, and the NeurIPS 2026 paper  
775 template. AIR’s `MetaChain ReAct` agent runs inside a Docker container; we configured the same  
776 Gemini 3.1 Pro model and iteration budgets as other systems. Of the 15 runs, 6 required re-runs due  
777 to Gemini API rate limits or sandbox crashes; 4 papers required manual  $\LaTeX$  fixes (trailing newlines,  
778 malformed `filecontents` environments, bibliography corruption).

## 779 **NeurIPS Paper Checklist**

### 780 **1. Claims**

781 Question: Do the main claims made in the abstract and introduction accurately reflect the  
782 paper’s contributions and scope?

783 Answer: [Yes]

784 Justification: Our claims are clearly justified by experimental results and discussion.

785 Guidelines:

- 786 • The answer [N/A] means that the abstract and introduction do not include the claims  
787 made in the paper.
- 788 • The abstract and/or introduction should clearly state the claims made, including the  
789 contributions made in the paper and important assumptions and limitations. A [No] or  
790 [N/A] answer to this question will not be perceived well by the reviewers.
- 791 • The claims made should match theoretical and experimental results, and reflect how  
792 much the results can be expected to generalize to other settings.
- 793 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
794 are not attained by the paper.

### 795 **2. Limitations**

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

797 Answer: [Yes]

798 Justification: We have discussed the detailed limitations in Appendix A.

799 Guidelines:

- 800 • The answer [N/A] means that the paper has no limitation while the answer [No] means  
801 that the paper has limitations, but those are not discussed in the paper.
- 802 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 803 • The paper should point out any strong assumptions and how robust the results are to  
804 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
805 model well-specification, asymptotic approximations only holding locally). The authors  
806 should reflect on how these assumptions might be violated in practice and what the  
807 implications would be.
- 808 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
809 only tested on a few datasets or with a few runs. In general, empirical results often  
810 depend on implicit assumptions, which should be articulated.
- 811 • The authors should reflect on the factors that influence the performance of the approach.  
812 For example, a facial recognition algorithm may perform poorly when image resolution  
813 is low or images are taken in low lighting. Or a speech-to-text system might not be  
814 used reliably to provide closed captions for online lectures because it fails to handle  
815 technical jargon.
- 816 • The authors should discuss the computational efficiency of the proposed algorithms  
817 and how they scale with dataset size.
- 818 • If applicable, the authors should discuss possible limitations of their approach to  
819 address problems of privacy and fairness.
- 820 • While the authors might fear that complete honesty about limitations might be used by  
821 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
822 limitations that aren’t acknowledged in the paper. The authors should use their best  
823 judgment and recognize that individual actions in favor of transparency play an impor-  
824 tant role in developing norms that preserve the integrity of the community. Reviewers  
825 will be specifically instructed to not penalize honesty concerning limitations.

### 826 **3. Theory assumptions and proofs**

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

829 Answer: [N/A]

830 Justification: Paper does not include theoretical results.

831 Guidelines:

- 832 • The answer [N/A] means that the paper does not include theoretical results.
- 833 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 834 referenced.
- 835 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 836 • The proofs can either appear in the main paper or the supplemental material, but if
- 837 they appear in the supplemental material, the authors are encouraged to provide a short
- 838 proof sketch to provide intuition.
- 839 • Inversely, any informal proof provided in the core of the paper should be complemented
- 840 by formal proofs provided in appendix or supplemental material.
- 841 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 842 4. Experimental result reproducibility

843 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

844 perimental results of the paper to the extent that it affects the main claims and/or conclusions

845 of the paper (regardless of whether the code and data are provided or not)?

846 Answer: [Yes]

847 Justification: We have provided necessary information for reproducibility as much as

848 possible in both main manuscript and appendix.

849 Guidelines:

- 850 • The answer [N/A] means that the paper does not include experiments.
- 851 • If the paper includes experiments, a [No] answer to this question will not be perceived
- 852 well by the reviewers: Making the paper reproducible is important, regardless of
- 853 whether the code and data are provided or not.
- 854 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 855 to make their results reproducible or verifiable.
- 856 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 857 For example, if the contribution is a novel architecture, describing the architecture fully
- 858 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 859 be necessary to either make it possible for others to replicate the model with the same
- 860 dataset, or provide access to the model. In general, releasing code and data is often
- 861 one good way to accomplish this, but reproducibility can also be provided via detailed
- 862 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 863 of a large language model), releasing of a model checkpoint, or other means that are
- 864 appropriate to the research performed.
- 865 • While NeurIPS does not require releasing code, the conference does require all submis-
- 866 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 867 nature of the contribution. For example
- 868 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
- 869 to reproduce that algorithm.
- 870 (b) If the contribution is primarily a new model architecture, the paper should describe
- 871 the architecture clearly and fully.
- 872 (c) If the contribution is a new model (e.g., a large language model), then there should
- 873 either be a way to access this model for reproducing the results or a way to reproduce
- 874 the model (e.g., with an open-source dataset or instructions for how to construct
- 875 the dataset).
- 876 (d) We recognize that reproducibility may be tricky in some cases, in which case
- 877 authors are welcome to describe the particular way they provide for reproducibility.
- 878 In the case of closed-source models, it may be that access to the model is limited in
- 879 some way (e.g., to registered users), but it should be possible for other researchers
- 880 to have some path to reproducing or verifying the results.

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

882 Question: Does the paper provide open access to the data and code, with sufficient instruc-

883 tions to faithfully reproduce the main experimental results, as described in supplemental

884 material?

885  
886  
887  
888  
889  
890  
891  
892  
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

Answer: [No]

Justification: We have provided detailed methodology section as well as appendix. We will consider code publication upon acceptance.

Guidelines:

- The answer [N/A] means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) 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. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 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) necessary to understand the results?

Answer: [Yes]

Justification: We have provided all the details in the appendix and the method section.

Guidelines:

- The answer [N/A] 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: We provide standard deviations of the experiments in the appendix.

Guidelines:

- The answer [N/A] 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, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- 937
- 938
- 939
- 940
- 941
- 942
- 943
- 944
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
  - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g., negative error rates).
  - If error bars are reported in tables or plots, the authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

945 **8. Experiments compute resources**

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

949 Answer: [Yes]

950 Justification: We provided detailed computational resources in the appendix.

951 Guidelines:

- 952
- 953
- 954
- 955
- 956
- 957
- 958
- 959
- The answer [N/A] 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.
  - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

960 **9. Code of ethics**

961 Question: Does the research conducted in the paper conform, in every respect, with the  
962 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

963 Answer: [Yes]

964 Justification: We followed the Code of Ethics.

965 Guidelines:

- 966
- 967
- 968
- 969
- 970
- 971
- The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
  - If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.
  - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

972 **10. Broader impacts**

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

975 Answer: [Yes]

976 Justification: We discussed potential societal impacts in the Appendix A.

977 Guidelines:

- 978
- 979
- 980
- 981
- 982
- 983
- 984
- The answer [N/A] means that there is no societal impact of the work performed.
  - If the authors answer [N/A] or [No], they should explain why their work has no societal impact or why the paper does not address societal impact.
  - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- 985
- 986
- 987
- 988
- 989
- 990
- 991
- 992
- 993
- 994
- 995
- 996
- 997
- 998
- 999
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
  - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 1000 11. Safeguards

1001 Question: Does the paper describe safeguards that have been put in place for responsible  
1002 release of data or models that have a high risk for misuse (e.g., pre-trained language models,  
1003 image generators, or scraped datasets)?

1004 Answer: [N/A]

1005 Justification: We have not released any data or models in this submission.

1006 Guidelines:

- 1007
- 1008
- 1009
- 1010
- 1011
- 1012
- 1013
- 1014
- 1015
- 1016
- The answer [N/A] means that the paper poses no such risks.
  - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
  - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
  - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 1017 12. Licenses for existing assets

1018 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
1019 the paper, properly credited and are the license and terms of use explicitly mentioned and  
1020 properly respected?

1021 Answer: [Yes]

1022 Justification: We have properly cited the assets used in the paper.

1023 Guidelines:

- 1024
- 1025
- 1026
- 1027
- 1028
- 1029
- 1030
- 1031
- 1032
- 1033
- 1034
- 1035
- 1036
- The answer [N/A] means that the paper does not use existing assets.
  - The authors should cite the original paper that produced the code package or dataset.
  - The authors should state which version of the asset is used and, if possible, include a URL.
  - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
  - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
  - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
  - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

1037 • If this information is not available online, the authors are encouraged to reach out to  
1038 the asset’s creators.

1039 **13. New assets**

1040 Question: Are new assets introduced in the paper well documented and is the documentation  
1041 provided alongside the assets?

1042 Answer: [N/A]

1043 Justification: This submission does not release new assets.

1044 Guidelines:

- 1045 • The answer [N/A] means that the paper does not release new assets.
- 1046 • Researchers should communicate the details of the dataset/code/model as part of their  
1047 submissions via structured templates. This includes details about training, license,  
1048 limitations, etc.
- 1049 • The paper should discuss whether and how consent was obtained from people whose  
1050 asset is used.
- 1051 • At submission time, remember to anonymize your assets (if applicable). You can either  
1052 create an anonymized URL or include an anonymized zip file.

1053 **14. Crowdsourcing and research with human subjects**

1054 Question: For crowdsourcing experiments and research with human subjects, does the paper  
1055 include the full text of instructions given to participants and screenshots, if applicable, as  
1056 well as details about compensation (if any)?

1057 Answer: [N/A].

1058 Justification: This paper does not involve crowdsourcing nor research with human subjects.

1059 Guidelines:

- 1060 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
1061 with human subjects.
- 1062 • Including this information in the supplemental material is fine, but if the main contribu-  
1063 tion of the paper involves human subjects, then as much detail as possible should be  
1064 included in the main paper.
- 1065 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
1066 or other labor should be paid at least the minimum wage in the country of the data  
1067 collector.

1068 **15. Institutional review board (IRB) approvals or equivalent for research with human  
1069 subjects**

1070 Question: Does the paper describe potential risks incurred by study participants, whether  
1071 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
1072 approvals (or an equivalent approval/review based on the requirements of your country or  
1073 institution) were obtained?

1074 Answer: [N/A]

1075 Justification: This paper does not involve crowdsourcing nor research with human subjects

1076 Guidelines:

- 1077 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
1078 with human subjects.
- 1079 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
1080 may be required for any human subjects research. If you obtained IRB approval, you  
1081 should clearly state this in the paper.
- 1082 • We recognize that the procedures for this may vary significantly between institutions  
1083 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
1084 guidelines for their institution.
- 1085 • For initial submissions, do not include any information that would break anonymity (if  
1086 applicable), such as the institution conducting the review.

1087 **16. Declaration of LLM usage**

1088 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
1089 non-standard component of the core methods in this research? Note that if the LLM is used  
1090 only for writing, editing, or formatting purposes and does *not* impact the core methodology,  
1091 scientific rigor, or originality of the research, declaration is not required.

1092 Answer: [\[Yes\]](#)

1093 Justification: We have provided the usages of the LLM in the Appendix.

1094 Guidelines:

- 1095 • The answer [\[N/A\]](#) means that the core method development in this research does not  
1096 involve LLMs as any important, original, or non-standard components.
- 1097 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not  
1098 be described.