
The Overfitting Crisis in LLM Workflows: Learning from Machine Learning’s Past Mistakes

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The rapid development of sophisticated Large Language Model (LLM) work-
2 flows—including agentic systems, multi-step reasoning pipelines, and tool-
3 integrated approaches—has led to impressive reported performance across various
4 benchmarks. However, we argue that the field is repeating a critical mistake from
5 early machine learning: reporting results on data that has been implicitly used for
6 training or optimization. The complexity of modern LLM workflows obscures the
7 fact that iterative prompt engineering, benchmark-driven development, and work-
8 flow refinement constitute a form of training on evaluation data. This position paper
9 draws parallels to historical overfitting practices in ML, documents how current
10 LLM development methodologies systematically conflate training and testing data,
11 and proposes best practices to address this growing methodological crisis before it
12 undermines the credibility of AI research, particularly in scientific applications.

13 1 Introduction

14 The rapid development of sophisticated Large Language Model (LLM) work-
15 flows—including agentic systems, multi-step reasoning pipelines, and tool-integrated approaches—has led to impressive
16 reported performance across various benchmarks [Bubeck et al., 2023, Wang et al., 2023]. However,
17 we argue that the field is repeating a critical mistake from early machine learning: reporting results
18 on data that has been implicitly used for training or optimization.

19 The complexity of modern LLM workflows obscures the fact that iterative prompt engineering,
20 benchmark-driven development, and workflow refinement constitute a form of training on evaluation
21 data [Min et al., 2022]. Contemporary LLM development increasingly involves sophisticated multi-
22 component systems where prompt optimization, tool selection, and architectural decisions are
23 iteratively refined based on performance feedback from target benchmarks [Schick and Schütze, 2020,
24 Qin et al., 2023]. This process mirrors the problematic practices that plagued early machine learning
25 research, where model selection and hyperparameter tuning were driven by test set performance.

26 Unlike traditional machine learning where the boundary between training and testing was relatively
27 clear, LLM workflow development operates in a gray area where the distinction between legitimate
28 system design and implicit training on evaluation data becomes blurred [Perez et al., 2021, Recht
29 et al., 2019]. The iterative nature of prompt engineering, the community-driven sharing of successful
30 strategies, and the progressive refinement of multi-agent architectures all contribute to a systematic
31 optimization process that effectively treats evaluation benchmarks as validation sets.

32 In early machine learning research, a common antipattern emerged that would later be recognized as
33 a fundamental threat to scientific validity. Researchers would select and tune models based on perfor-
34 mance on available datasets, with hyperparameter optimization, feature engineering, and architecture
35 choices driven by test set performance [Dwork et al., 2015, Blum and Hardt, 2015]. Final results
36 were reported on the same datasets used for development, leading to real-world deployment showing

37 significant performance degradation—a phenomenon that became known as the "reproducibility
38 crisis" in machine learning [Henderson et al., 2018, Lipton and Steinhardt, 2019].

39 The machine learning community's recognition of this methodological error represents one of the
40 most important paradigm shifts in computational research methodology. The "test set" had effectively
41 become part of the training process through repeated evaluation and optimization cycles, violating
42 the fundamental assumption of independent evaluation [Ioannidis, 2005]. This realization catalyzed
43 the establishment of rigorous protocols that became foundational to credible ML research: clear
44 separation of training, validation, and test sets; test sets remaining untouched until final evaluation;
45 cross-validation and holdout validation for development; and independent evaluation on truly unseen
46 data [Kohavi, 1995, Varma and Simon, 2006].

47 These methodological advances transformed machine learning from a field prone to optimistic bias
48 into one with robust evaluation standards [Raschka, 2018]. However, the emergence of complex
49 LLM workflows has created new opportunities for the same fundamental errors to re-emerge in more
50 sophisticated forms. This position paper draws parallels to historical overfitting practices in ML,
51 documents how current LLM development methodologies systematically conflate training and testing
52 data, and proposes best practices to address this growing methodological crisis before it undermines
53 the credibility of AI research, particularly in scientific applications.

54 **2 Related Work**

55 **2.1 Historical Overfitting in Machine Learning**

56 The machine learning community has long recognized the dangers of overfitting to evaluation data.
57 Early ML research suffered from researchers repeatedly testing models on the same evaluation sets,
58 leading to inflated performance estimates that failed to generalize [Dwork et al., 2015]. This led to
59 the establishment of rigorous evaluation protocols including proper train/validation/test splits and
60 holdout methodologies.

61 The development of comprehensive benchmarking frameworks helped standardize evaluation prac-
62 tices [Olson et al., 2017]. However, as noted by Cohen-Inger et al. [2025], benchmark datasets
63 themselves can become overoptimized targets when the entire research community focuses on the
64 same evaluation sets over extended periods.

65 **2.2 Data Contamination in Large Language Models**

66 Recent work has extensively documented data contamination issues in LLM training and evaluation.
67 Golchin and Surdeanu [2023] and Shi et al. [2024] demonstrate that simple decontamination meth-
68 ods like n-gram matching are insufficient, as paraphrasing and translation can easily bypass these
69 measures.

70 The LessLeak-Bench study [Zhou et al., 2025] provided comprehensive analysis across 83 software
71 engineering benchmarks, finding contamination across nearly all tested models. This systematic
72 contamination undermines reported performance improvements in domains where sophisticated LLM
73 workflows are increasingly deployed.

74 **2.3 Evaluation Methodology in ML Systems**

75 The establishment of rigorous evaluation protocols emerged as the ML community's primary defense
76 against overfitting. The foundational train/validation/test split methodology became standard practice,
77 with cross-validation techniques providing robust model selection frameworks [Stone, 1974, Kohavi,
78 1995]. These approaches ensured that model selection decisions were made independently of final
79 performance evaluation, preventing the test set from becoming an implicit part of the training process
80 [Hastie et al., 2009, Bishop, 2006].

81 Modern ML systems have extended these principles to encompass continuous evaluation and mon-
82 itoring throughout the system lifecycle. However, the fundamental principle remains unchanged:
83 evaluation data must remain independent of optimization processes to ensure valid performance
84 estimates.

85 **3 The Current Crisis in LLM Workflows**

86 **3.1 The Implicit Training Problem**

87 Modern LLM workflows involve extensive optimization processes that constitute implicit training on
88 evaluation data. Teams spend weeks iterating on prompts, testing variants against target benchmarks,
89 and selecting the best-performing approaches [Wei et al., 2022]. Each iteration uses the benchmark
90 as a validation signal.

91 Multi-agent systems, RAG pipelines, and tool usage patterns are designed and refined based on
92 performance on specific tasks and datasets. The architecture itself becomes optimized for known
93 evaluation criteria. Research teams explicitly target improvements on established benchmarks, using
94 them as development objectives rather than independent evaluation metrics.

95 **3.2 Complexity as Camouflage**

96 The sophistication of modern LLM workflows obscures the overfitting problem. Complex chains of
97 thought, tree search, and multi-agent interactions create an illusion of generalization, when in fact the
98 entire pipeline has been optimized for specific evaluation patterns [Yao et al., 2023].

99 Systems that use calculators, code interpreters, and web search appear more general, but tool
100 selection, usage patterns, and integration strategies are typically optimized against known benchmarks.
101 Workflows that adapt their strategies based on input characteristics seem more robust, but these
102 adaptation mechanisms are usually developed and tuned using the same evaluation data they will
103 later be tested on.

104 **3.3 Benchmark Contamination and Spillover Effects**

105 The problem is exacerbated by several forms of contamination. Many benchmarks overlap with or
106 derive from data sources used in LLM pre-training, creating subtle forms of data leakage [Shi et al.,
107 2024]. Recent work has shown that simple paraphrasing can bypass decontamination measures, and
108 that models can achieve artificially high performance when such variations are not eliminated [Shi
109 et al., 2024].

110 Solutions and techniques developed for one benchmark quickly propagate to others, creating implicit
111 optimization across multiple evaluation sets [Rogers et al., 2021, Dodge et al., 2021]. As models
112 improve on existing benchmarks, new versions are created that often share similar patterns and
113 evaluation criteria, perpetuating the contamination cycle.

114 **4 Detailed Case Studies: Workflow Overfitting and Generalisability
115 Challenges in Practice**

116 **4.1 ReAct: Reasoning and Acting Workflows**

117 The ReAct framework [Yao et al., 2022] exemplifies how sophisticated LLM workflows can be
118 systematically overfitted to evaluation benchmarks. ReAct combines reasoning traces with action-
119 taking capabilities, enabling LLMs to interact with external tools while maintaining step-by-step
120 reasoning.

121 **Development Process Analysis:** The ReAct paper reports evaluation on four benchmarks: HotpotQA
122 for question answering, Fever for fact verification, ALFWorld for text-based games, and WebShop for
123 web navigation. The development methodology involved iterative prompt engineering and workflow
124 refinement specifically targeting these benchmarks.

125 **Overfitting Evidence:** The ReAct development process exhibits several characteristics indicative of
126 benchmark optimization:

- 127 • Systematic testing of different prompting strategies against the target benchmarks
128 • Iterative refinement of the reasoning-action interleaving based on benchmark performance
129 • Optimization of tool integration strategies for benchmark-specific requirements

- 130 • Fine-tuning of termination criteria based on evaluation results

131 **Performance Claims:** The paper reports substantial improvements over baselines: ReAct achieves
132 27.4% success on HotpotQA compared to 20.6% for chain-of-thought prompting, representing a
133 33% relative improvement [Yao et al., 2022]. However, these gains emerged through extensive
134 optimization against these specific evaluation sets.

135 **4.2 Code Generation Workflows**

136 Code generation tasks present particularly clear examples of workflow overfitting. Benchmarks like
137 HumanEval [Chen et al., 2021] and MBPP [Austin et al., 2021] have become central targets for LLM
138 development, with teams iteratively refining their approaches against these specific problems.

139 **Development Process Analysis:** Modern coding workflows undergo extensive optimization cycles:

- 140 • Tool integration strategies refined based on benchmark performance
141 • Error handling approaches optimized for common benchmark failure modes
142 • Multi-step reasoning patterns tuned to handle benchmark-specific problem structures
143 • Code execution and debugging workflows designed around benchmark evaluation criteria

144 **Systematic Contamination Evidence:** The LessLeak-Bench study [Zhou et al., 2025] provides
145 compelling evidence of widespread contamination across software engineering benchmarks. The
146 study analyzed 83 benchmarks and found average leakage ratios of 4.8%, 2.8%, and 0.7% for Python,
147 Java, and C/C++ benchmarks respectively. However, specific benchmarks showed much higher
148 contamination rates, with QuixBugs exhibiting 100% leakage and BigCloneBench showing 55.7%
149 leakage.

150 **Performance Impact:** The study demonstrates that data leakage has substantial impact on LLM
151 evaluation, with contaminated models showing inflated performance that does not generalize to truly
152 novel programming challenges.

153 **4.3 Autonomous Agent Systems**

154 Autonomous agent frameworks like AutoGPT represent complex multi-component systems that are
155 particularly susceptible to overfitting due to their iterative development processes and community-
156 driven optimization.

157 **Development Challenges:** Agent systems face unique evaluation challenges:

- 158 • Benchmarks developed concurrently with the systems they evaluate
159 • Community-driven optimization leading to distributed overfitting effects
160 • Performance claims based on evaluation sets that guided development decisions
161 • Informal evaluation criteria that evolve based on system capabilities

162 **Generalization Gaps:** Despite reported success on development benchmarks, deployed agent systems
163 frequently exhibit significant performance degradation when encountering novel scenarios that differ
164 from their optimization targets [Zheng et al., 2023, Garg et al., 2025]. This pattern suggests that
165 benchmark performance may reflect sophisticated pattern matching rather than genuine autonomous
166 reasoning capabilities.

167 **4.4 Chemistry and Materials Science**

168 Scientific applications of LLM workflows present particularly high stakes for the overfitting problem.
169 Recent work in chemistry has developed sophisticated benchmarks like ChemBench [Mirza et al.,
170 2024], containing over 2,700 question-answer pairs designed to evaluate chemical knowledge and
171 reasoning. While these benchmarks represent important advances in evaluation methodology, they
172 also create new targets for optimization.

173 **Domain-Specific Overfitting Risks:** Chemical reasoning workflows often involve:

- 174 • Multi-step synthesis planning optimized against known reaction databases
175 • Property prediction systems trained and validated on established chemical datasets
176 • Literature analysis tools refined using benchmark chemical papers and abstracts
177 • Safety assessment frameworks tuned against standardized hazard classification systems

178 **Reproducibility Standards in Chemistry:** The chemistry community has established rigorous
179 standards for experimental reproducibility, but these standards have not yet been adapted for AI-
180 assisted chemical research. Unlike traditional chemistry experiments, where failed replications are
181 clearly identifiable, overfitted AI systems may produce plausible but incorrect results that are difficult
182 to detect without expert domain knowledge.

183 **4.5 Physics and Engineering Applications**

184 Physics applications face similar challenges, with benchmarks targeting graduate-level physics prob-
185 lems. Engineering applications often require reasoning across multiple domains simultaneously, but
186 systems optimized against such benchmarks may develop sophisticated pattern-matching capabilities
187 that fail when confronted with truly novel interdisciplinary problems [Zhang et al., 2025b,a].

188 **Multi-Domain Scientific Reasoning:** Scientific applications often require reasoning across multiple
189 domains simultaneously [Cui et al., 2025], creating additional opportunities for overfitting to develop
190 on multi-faceted benchmark tasks.

191 **4.6 Biomedical and Healthcare Applications**

192 Healthcare applications represent the highest-stakes domain for AI system reliability. Biomedical
193 LLM workflows increasingly target specialized benchmarks for clinical reasoning, drug discovery,
194 and diagnostic assistance [Zhu et al., 2025, Tang et al., 2025].

195 **Critical Safety Implications:** In healthcare contexts, overfitted performance can have severe conse-
196 quences:

- 197 • Diagnostic systems optimized against medical benchmark datasets may miss novel disease
198 presentations
199 • Drug discovery workflows trained on established chemical databases may fail to identify
200 truly innovative therapeutic approaches
201 • Clinical decision support systems refined using benchmark cases may provide inappropriate
202 recommendations for edge cases

203 **Regulatory and Compliance Challenges:** Healthcare AI systems must meet strict regulatory
204 standards for safety and efficacy. However, current evaluation practices may not adequately dis-
205 tinguish between genuine clinical reasoning and sophisticated pattern matching against medical
206 benchmarks [Bedi et al., 2025, Mehandru et al., 2025].

207 **5 Comprehensive Framework for Addressing Workflow Overfitting**

208 Building on established ML evaluation best practices, we propose a comprehensive approach to
209 address workflow overfitting that extends beyond traditional model evaluation to encompass the entire
210 development ecosystem.

211 **5.1 Rigorous Data Separation Methodologies**

212 **Fundamental Train/Validation/Test Separation:** LLM workflow development must adopt the foun-
213 dational ML principle of strict data separation. Training data is used for initial system development,
214 validation data guides iterative refinement and hyperparameter optimization, and test data remains
215 completely isolated until final evaluation. This separation must be maintained throughout the entire
216 workflow development process, including prompt engineering, tool integration, and architectural
217 decisions.

218 **Temporal Isolation Protocols:** Evaluation benchmarks must use data that postdates workflow
219 development completion. This temporal separation prevents both direct benchmark exposure and
220 indirect contamination through community knowledge dissemination. The temporal gap must account
221 for publication delays and community adoption cycles, ensuring no development decisions can be
222 influenced by evaluation content.

223 **Cross-Domain Transfer Evaluation:** Systems developed in one domain should be evaluated on
224 structurally analogous tasks in different domains. For instance, reasoning workflows optimized on
225 historical datasets should be tested on scientific reasoning problems with similar logical structures
226 but distinct knowledge bases. This tests genuine transferable capabilities rather than domain-specific
227 pattern recognition.

228 **Hierarchical Holdout Architecture:** Implement nested isolation strategies at multiple levels:
229 component-level evaluation using standard ML holdout practices, workflow-level assessment on
230 integration benchmarks isolated from development, and system-level evaluation on deployment
231 scenarios that differ systematically from development contexts.

232 **5.2 Methodology Transparency Requirements**

233 Current LLM workflow development lacks the systematic documentation that enables contamination
234 detection and replication. Unlike traditional ML where training procedures are explicitly documented,
235 workflow optimization often occurs through informal iteration cycles.

236 **Development History Documentation:** Comprehensive logging must capture all evaluation data
237 interactions, including datasets accessed during prompt engineering, optimization iterations against
238 benchmarks, architectural decisions guided by performance feedback, and community knowledge
239 sources consulted.

240 **Contamination Auditing:** Automated analysis must detect potential overlap between development
241 resources and evaluation sets, including computational analysis of data intersection, documentation
242 of influential published work, and assessment of community knowledge transfer.

243 **Preregistration Protocols:** Adapting practices from experimental psychology and medical research,
244 teams must declare evaluation benchmarks before development, specify success criteria in advance,
245 and commit to reporting results regardless of outcome.

246 **5.3 Institutional and Ecosystem-Level Interventions**

247 **Benchmark Lifecycle Management:** Establish systematic protocols for benchmark retirement
248 and replacement, including monitoring community, implementing retirement triggers based on
249 performance saturation, and developing principled approaches for creating replacement benchmarks.

250 **Independent Evaluation Infrastructure:** Create community-managed evaluation services main-
251 taining truly independent benchmarks inaccessible to development teams until final evaluation. This
252 requires sustained funding and governance structures resistant to commercial and academic pressures.

253 **Research Incentive Realignment:** Conference and journal policies must require comprehensive
254 methodology transparency. Recognition systems should value negative results and replication studies
255 equally with novel contributions. Funding agencies must prioritize evaluation methodology alongside
256 technical innovation.

257 **6 Why This Matters for Scientific Applications**

258 The overfitting problem is particularly concerning for scientific applications. Scientific research
259 demands reproducible results, but if LLM workflows are optimized against the same benchmarks
260 they're evaluated on, reported performance may not generalize to real scientific problems.

261 Overfitted workflows may perform well on known benchmarks but fail on novel scientific challenges,
262 leading to misplaced confidence in AI capabilities for scientific discovery. If reported performance is
263 inflated due to overfitting, research funding and effort may be misdirected toward approaches that
264 don't actually advance scientific capability.

265 The stakes are particularly high in scientific domains where accuracy and reliability are paramount,
266 and where the cost of false confidence can impede genuine scientific progress.

267 7 Limitations

268 This position paper has several important limitations that should be acknowledged when interpreting
269 our arguments and recommendations.

270 **Lack of Empirical Validation:** Our analysis relies primarily on theoretical reasoning and examination
271 of existing literature rather than systematic empirical studies. While we cite evidence from studies
272 like LessLeak-Bench [Zhou et al., 2025], we have not conducted controlled experiments to directly
273 measure workflow overfitting or validate our proposed solutions.

274 **Limited Access to Proprietary Systems:** Our case studies focus on publicly documented systems
275 and open benchmarks. Many commercial LLM workflows involve proprietary development processes
276 that are not accessible for analysis, limiting our ability to assess the full scope of the problem across
277 the industry.

278 **Implementation Feasibility:** While we propose comprehensive solutions, we acknowledge sig-
279 nificant practical barriers including costs and coordination challenges. The feasibility of these
280 recommendations remains untested.

281 **Boundary Definition:** The distinction between legitimate iterative development and problematic
282 overfitting is not always clear-cut. Our framework does not provide precise operational definitions
283 for this boundary, which could lead to either overly restrictive or insufficient practices.

284 8 Discussion and Implications

285 8.1 Systemic Nature of Workflow Overfitting

286 Unlike traditional ML overfitting, which typically affected individual models, LLM workflow overfit-
287 ting operates at an unprecedented scale with systemic consequences. When the research community
288 collectively optimizes against the same benchmarks, the cumulative effect creates systematic bias that
289 permeates entire research directions. This represents "ecosystem-level overfitting" where distributed
290 optimization across hundreds of research teams amplifies the problem beyond classical overfitting
291 scenarios.

292 The interconnected nature of modern AI research exacerbates this through rapid technique propagation
293 via preprints and code repositories. Competitive pressure to achieve state-of-the-art results on
294 established leaderboards incentivizes optimization for known evaluation sets rather than developing
295 genuinely novel capabilities, creating a feedback loop where benchmark performance becomes
296 divorced from real-world utility.

297 8.2 Epistemological and Economic Implications

298 Current practices fundamentally conflate innovation with evaluation, undermining the epistemological
299 foundations of AI research. True scientific evaluation requires independence from development—a
300 principle fundamental to empirical inquiry. When evaluation data contaminates development deci-
301 sions, resulting performance metrics cannot serve as valid evidence for genuine capability advances.

302 This confusion has profound economic implications as organizations increasingly rely on benchmark
303 performance to guide technology adoption and investment decisions. Overfitted results can lead to
304 substantial resource misallocation, with companies investing heavily in systems that perform well on
305 benchmarks but fail in production environments.

306 8.3 High-Stakes Domain Risks

307 As LLM workflows become more sophisticated and deployed in critical applications—from drug
308 discovery to climate modeling—the cost of overfitted performance becomes exponentially higher.
309 Scientific discovery relies on generating novel insights and identifying previously unknown patterns,

310 capabilities fundamentally undermined when systems optimize primarily for existing datasets that
311 may not capture natural phenomena's full complexity.
312 The risk is particularly severe where ground truth is difficult to establish or system failures may not
313 be immediately apparent. In chemistry, for example, AI systems excelling at predicting molecular
314 properties on benchmark datasets might fail catastrophically with novel compounds, potentially
315 delaying scientific breakthroughs or causing harmful outcomes if trusted beyond actual capabilities.
316 The scientific and commercial communities cannot afford to repeat machine learning's early method-
317 ological mistakes, particularly when applications involve high-stakes scientific discovery and decision-
318 making where system failures can have serious real-world consequences.

319 **9 Conclusion**

320 The LLM research community stands at a critical juncture. The sophistication of modern workflows
321 has obscured a fundamental methodological problem: we are systematically reporting results on data
322 that has been used for optimization. This mirrors the overfitting crisis that plagued early machine
323 learning research.
324 Unlike traditional ML, where overfitting affected individual models, LLM workflow overfitting can
325 contaminate entire research directions and mislead the community about actual capabilities. This is
326 particularly concerning for scientific applications, where accuracy and reliability are paramount.
327 The solution requires collective action: establishing strict data separation protocols, demanding
328 transparency in development processes, and creating truly independent evaluation resources. The ML
329 community learned these lessons decades ago. The LLM community must learn them now, before
330 the credibility of AI research is further undermined.
331 We call on the research community to adopt rigorous evaluation standards that separate development
332 from testing, just as was necessary in traditional ML. Only through such methodological rigor can
333 we ensure that reported advances represent genuine progress rather than sophisticated forms of
334 overfitting.

335 **References**

- 336 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,
337 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language
338 models. *arXiv preprint arXiv:2108.07732*, 2021.
- 339 Suhana Bedi, Yixing Jiang, Philip Chung, Sanmi Koyejo, and Nigam Shah. Fidelity of medical
340 reasoning in large language models. *JAMA Network Open*, 8(8):e2526021–e2526021, 2025.
- 341 Christopher M Bishop. *Pattern recognition and machine learning*. Springer, 2006.
- 342 Avrim Blum and Moritz Hardt. The ladder: A reliable leaderboard for machine learning competitions.
343 *International conference on machine learning*, pages 1006–1014, 2015.
- 344 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar,
345 Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence:
346 Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- 347 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared
348 Kaplan, Harri Edwards, Yura Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large
349 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 350 Nurit Cohen-Inger, Yehonatan Elisha, Bracha Shapira, Lior Rokach, and Seffi Cohen. Forget what
351 you know about llms evaluations—llms are like a chameleon. *arXiv preprint arXiv:2502.07445*,
352 2025.
- 353 Hao Cui, Zahra Shamsi, Gowoon Cheon, Xuejian Ma, Shutong Li, Maria Tikhonovskaya, Peter
354 Norgaard, Nayantara Mudur, Martyna Plomecka, Paul Raccuglia, et al. Curie: Evaluating llms on
355 multitask scientific long context understanding and reasoning. *arXiv preprint arXiv:2503.13517*,
356 2025.

- 357 Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld,
358 Margaret Mitchell, and Matt Gardner. Documenting large webtext corpora: A case study on the
359 colossal clean crawled corpus. In *Proceedings of the 2021 Conference on Empirical Methods in*
360 *Natural Language Processing*, pages 1286–1305, 2021.
- 361 Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Roth.
362 The reusable holdout: Preserving validity in adaptive data analysis. *Science*, 349(6248):636–638,
363 2015.
- 364 Divyansh Garg, Shaun VanWeelden, Diego Caples, Andis Draguns, Nikil Ravi, Pranav Putta, Naman
365 Garg, Tomas Abraham, Michael Lara, Federico Lopez, et al. Real: Benchmarking autonomous
366 agents on deterministic simulations of real websites. *arXiv preprint arXiv:2504.11543*, 2025.
- 367 Shahriar Golchin and Mihai Surdeanu. Time travel in llms: Tracing data contamination in large
368 language models. *arXiv preprint arXiv:2308.08493*, 2023.
- 369 Trevor Hastie, Robert Tibshirani, and Jerome H Friedman. *The elements of statistical learning: data*
370 *mining, inference, and prediction*, volume 2. Springer, 2009.
- 371 Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger.
372 Deep reinforcement learning that matters. *Proceedings of the AAAI conference on artificial*
373 *intelligence*, 32(1), 2018.
- 374 John PA Ioannidis. Why most published research findings are false. *PLoS medicine*, 2(8):e124, 2005.
- 375 Ron Kohavi. A study of cross-validation and bootstrap for accuracy estimation and model selection.
376 In *Ijcai*, volume 14, pages 1137–1145, 1995.
- 377 Zachary C Lipton and Jacob Steinhardt. Troubling trends in machine learning scholarship: Some ml
378 papers suffer from flaws that could mislead the public and stymie future research. *Communications*
379 *of the ACM*, 62(6):45–53, 2019.
- 380 Nikita Mehandru, Niloufar Golchini, David Bamman, Travis Zack, Melanie F Molina, and Ahmed
381 Alaa. Er-reason: A benchmark dataset for llm-based clinical reasoning in the emergency room.
382 *arXiv preprint arXiv:2505.22919*, 2025.
- 383 Sewon Min, Xinxixi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke
384 Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv*
385 *preprint arXiv:2202.12837*, 2022.
- 386 Adrian Mirza, Nawaf Alampara, Sreekanth Kunchapu, Martíño Ríos-García, Benedict Emoekabu,
387 Aswanth Krishnan, Tanya Gupta, Mara Schilling-Wilhelmi, Macjonathan Okereke, Anagha Aneesh,
388 et al. Are large language models superhuman chemists? *arXiv preprint arXiv:2404.01475*, 2024.
- 389 Randal S Olson, William La Cava, Zairah Mustahsan, Akshay Varik, and Jason H Moore. Pmlb: a
390 large benchmark suite for machine learning evaluation and comparison. *BioData mining*, 10(1):
391 1–13, 2017.
- 392 Ethan Perez, Douwe Kiela, and Kyunghyun Cho. True few-shot learning with language models.
393 *Advances in neural information processing systems*, 34:11054–11070, 2021.
- 394 Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru
395 Tang, Bill Qian, et al. Tool learning with foundation models. *arXiv preprint arXiv:2304.08354*,
396 2023.
- 397 Sebastian Raschka. Model evaluation, model selection, and algorithm selection in machine learning.
398 *arXiv preprint arXiv:1811.12808*, 2018.
- 399 Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers
400 generalize to imagenet? *International conference on machine learning*, pages 5389–5400, 2019.
- 401 Anna Rogers, Timothy Baldwin, and Kobi Leins. Changing the world by changing the data. In
402 *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the*
403 *11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*,
404 pages 2182–2194, 2021.

- 405 Timo Schick and Hinrich Schütze. Exploiting cloze questions for few-shot text classification and
406 natural language inference. *arXiv preprint arXiv:2001.07676*, 2020.
- 407 Shuo Shi, Xiangyu Ku, Xiang Yue, Tianle Li, Yunhua Li, and William Yang Wang. Rethinking
408 benchmark and contamination for language models with rephrased samples. *arXiv preprint*
409 *arXiv:2311.04850*, 2024.
- 410 Mervyn Stone. Cross-validatory choice and assessment of statistical predictions. *Journal of the Royal*
411 *Statistical Society: Series B (Methodological)*, 36(2):111–133, 1974.
- 412 Xiangru Tang, Daniel Shao, Jiwoong Sohn, Jiapeng Chen, Jiayi Zhang, Jinyu Xiang, Fang Wu, Yilun
413 Zhao, Chenglin Wu, Wenqi Shi, et al. Medagentsbench: Benchmarking thinking models and agent
414 frameworks for complex medical reasoning. *arXiv preprint arXiv:2503.07459*, 2025.
- 415 Sudhir Varma and Richard Simon. Bias in error estimation when using cross-validation for model
416 selection. *BMC bioinformatics*, 7(1):1–8, 2006.
- 417 Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai
418 Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents.
419 *arXiv preprint arXiv:2308.11432*, 2023.
- 420 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, and
421 Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*
422 *Neural Information Processing Systems*, 35:24824–24837, 2022.
- 423 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
424 React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*,
425 2022.
- 426 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik
427 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances*
428 *in Neural Information Processing Systems*, 36, 2023.
- 429 Xinyu Zhang, Yuxuan Dong, Yanrui Wu, Jiaxing Huang, Chengyou Jia, Basura Fernando, Mike Zheng
430 Shou, Lingling Zhang, and Jun Liu. Physreason: A comprehensive benchmark towards physics-
431 based reasoning. *arXiv preprint arXiv:2502.12054*, 2025a.
- 432 Yiming Zhang, Yingfan Ma, Yanmei Gu, Zhengkai Yang, Yihong Zhuang, Feng Wang, Zenan
433 Huang, Yuanyuan Wang, Chao Huang, Bowen Song, et al. Abench-physics: Benchmarking
434 physical reasoning in llms via high-difficulty and dynamic physics problems. *arXiv preprint*
435 *arXiv:2507.04766*, 2025b.
- 436 Xinyue Zheng, Haowei Lin, Kaichen He, Zihao Wang, Zilong Zheng, QIANG FU, Haobo Fu, and
437 Yitao Liang. Towards evaluating generalist agents: An automated benchmark in open world. 2023.
- 438 Xin Zhou, Kisub Li, Yunbo Lyu, Ming Wen, Beijun Wang, Zhou Yang, Zhiqiang Wang, Xiaodong
439 Luo, Li Li, and Yang Li. Lessleak-bench: A first investigation of data leakage in llms across 83
440 software engineering benchmarks. *arXiv preprint arXiv:2502.06215*, 2025.
- 441 Yakun Zhu, Zhongzhen Huang, Linjie Mu, Yutong Huang, Wei Nie, Jiaji Liu, Shaoting Zhang,
442 Pengfei Liu, and Xiaofan Zhang. Diagnosisarena: Benchmarking diagnostic reasoning for large
443 language models. *arXiv preprint arXiv:2505.14107*, 2025.

444 **Agents4Science AI Involvement Checklist**

- 445 1. **Hypothesis development:** The research topic and hypothesis that LLM workflows suffer
446 from systematic overfitting similar to early ML practices was jointly developed, with the AI
447 agent providing the detailed conceptual framework and historical parallels.

448 Answer: **[B] Mostly human, assisted by AI**

449 Explanation: The human co-author identified the core problem and made the connection
450 to historical ML overfitting. The AI agent developed the detailed hypothesis, literature
451 connections, and theoretical framework for understanding the problem.

- 452 2. **Experimental design and implementation:** This is a position paper that does not include
453 empirical experiments. The analysis relies on synthesizing existing literature and making
454 methodological arguments.

455 Answer: **[NA] Not Applicable**

456 Explanation: No experiments were conducted. The paper makes its argument through
457 literature analysis, theoretical reasoning, and methodological critique of current practices.

- 458 3. **Analysis of data and interpretation of results:** The paper analyzes existing literature and
459 research practices rather than experimental data. The AI agent conducted literature search
460 and synthesis.

461 Answer: **[C] Mostly AI, assisted by human**

462 Explanation: The AI agent conducted web searches for relevant literature, identified key
463 papers on data contamination and benchmark leakage, and synthesized findings to support
464 the core arguments. Human provided guidance on focus areas.

- 465 4. **Writing:** The AI agent was responsible for the majority of the writing, including structure,
466 arguments, and academic formatting.

467 Answer: **[C] Mostly AI, assisted by human**

468 Explanation: The AI agent wrote the complete paper including abstract, introduction,
469 literature review, arguments, and conclusions. The human co-author provided feedback,
470 guidance on emphasis areas, and editorial input.

- 471 5. **Observed AI Limitations:** The AI agent required guidance on specific examples and
472 emphasis areas, and needed human oversight to ensure arguments remained grounded and
473 credible.

474 Description: The AI agent occasionally needed direction on which aspects of the overfitting
475 problem to emphasize most strongly. The human co-author's role was crucial in keeping the
476 argument focused and ensuring it addressed the most important methodological concerns in
477 the field.

478 **Agents4Science Paper Checklist**

479 **1. Claims**

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

482 Answer: [Yes]

483 Justification: The abstract and introduction clearly state this is a position paper arguing that
484 LLM workflows suffer from systematic overfitting, with claims appropriately scoped to
485 methodological critique rather than empirical findings.

486 **2. Limitations**

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

488 Answer: [Yes]

489 Justification: Section 7 discusses the scope of the argument and acknowledges this is a
490 methodological position paper without empirical validation of proposed solutions.

491 **3. Theory assumptions and proofs**

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

494 Answer: [NA]

495 Justification: This is a position paper that makes methodological arguments rather than
496 theoretical claims requiring formal proofs.

497 **4. Experimental result reproducibility**

498 Question: Does the paper fully disclose all the information needed to reproduce the main
499 experimental results?

500 Answer: [NA]

501 Justification: This paper does not include experiments. It is a methodological position paper
502 based on literature analysis and reasoned argument.

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

504 Question: Does the paper provide open access to the data and code?

505 Answer: [NA]

506 Justification: No experiments or code are involved in this position paper.

507 **6. Experimental setting/details**

508 Question: Does the paper specify all training and test details necessary to understand the
509 results?

510 Answer: [NA]

511 Justification: This paper does not include experiments.

512 **7. Experiment statistical significance**

513 Question: Does the paper report error bars or statistical significance information?

514 Answer: [NA]

515 Justification: No experiments are conducted in this position paper.

516 **8. Experiments compute resources**

517 Question: Does the paper provide sufficient information on computer resources needed?

518 Answer: [NA]

519 Justification: No experiments requiring computational resources were conducted.

520 **9. Code of ethics**

521 Question: Does the research conform with the Agents4Science Code of Ethics?

522 Answer: [Yes]

523 Justification: This methodological critique aims to improve research practices and scientific
524 integrity, fully conforming with ethical research standards.

525 **10. Broader impacts**

526 Question: Does the paper discuss potential positive and negative societal impacts?

527 Answer: [Yes]

528 Justification: Throughout the paper, we discuss both the positive impact of addressing
529 workflow overfitting (improving research credibility, advancing genuine scientific capa-
530 bilities) and the negative consequences of failing to address it (misplaced confidence in
531 AI systems, potential safety risks in high-stakes domains like healthcare, misdirection of
532 research resources).