

# AGENTIC UNCERTAINTY REVEALS AGENTIC OVER-CONFIDENCE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Can AI agents predict whether they will succeed at a task? We study **agentic uncertainty** by eliciting success probability estimates before, during, and after task execution. All results exhibit **agentic overconfidence**: some agents that succeed only 22% of the time predict 77% success. Counterintuitively, pre-execution assessment with strictly less information achieves better discrimination than standard post-execution review. Adversarial prompting reframing assessment as bug-finding achieves the best calibration.

## 1 INTRODUCTION

A software engineer needs to fix an auth service error. Before delegating to an AI coding agent, she asks: what are the chances this succeeds?

### Pre-Exec.

**P(success): 72%**

The issue is clear and well-defined, the error message points directly to the problem, and the fix follows existing patterns in the codebase.

72% confidence before any code is written. As the coding agent works, she asks another agent to monitor progress:

### Mid-Exec.

**P(success): 78%**

The agent has correctly diagnosed the problem and knows exactly what code to add. The probability of success is high.

Confidence *rises* to 78%. The patch is now complete. She fires off a review agent:

### Post-Exec.

**P(success): 92%**

The patch is a correct and complete fix. It's a minimal, focused change that adds the missing interface method.

Too optimistic. Let's spawn an adversarial agent.

### Adv. Post-Exec.

**P(success): 85%**

Minor concerns don't affect the main use case... the patch correctly resolves the reported issue.

Still 85%. All four agents confidently predict success.

But the patch fails! And this **agentic overconfidence** is systematic (Figure 1). For example, GPT-5.2-Codex-based post-execution agents predict 73% success against a true rate of 35% averaged over 100 SWE-Bench-Pro (Deng et al., 2025) tasks.

This matters because the scope of autonomous work is expanding rapidly. The effective length of tasks that AI agents complete has doubled every 7 months for six years (METR, 2025). As we

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

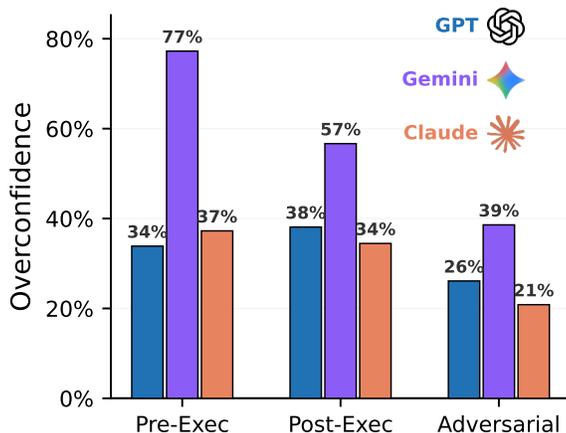


Figure 1: **Agentic overconfidence.** Overconfidence (predicted minus true success rate) across three strategies. True rates: GPT-5.2 Codex 35%, Gemini-3-Pro 22%, Opus 4.5 27%. All agents systematically overestimate their success.

increasingly delegate complex workflows to agents (Appel et al., 2025), we must develop scalable oversight protocols (Bowman et al., 2022).

In this work, we elicit *agentic uncertainty* at three points in a coding agent’s lifecycle: pre-, mid-, and post-execution. Each corresponds to a different oversight question: Can agents predict failure before committing resources? Can they recognize failure as it unfolds? Can they verify their own work? Importantly, we use the same underlying model for both the coding agent that produces patches and the uncertainty agent that assesses them, isolating the effect of information access from differences in model capability.

Our experiments on 100 SWE-bench Pro tasks across three frontier models (GPT-5.2-Codex, Gemini-3-Pro, Claude Opus 4.5) reveal several striking findings:

- **Pervasive overconfidence.** Post-execution agents can predict 73% success on average against a 35% base rate (GPT), with similar gaps across all models.
- **More context, uncalibrated doubt.** Mid-execution agents develop “cold feet”: confidence decreases as they observe their partial work, but this doubt is uninformative, occurring equally for successes and failures.
- **Adversarial framing helps.** Prompting agents to “find bugs” rather than “verify correctness” reduces overconfidence by up to 15 pp and achieves the best calibration across all models.

## 2 METHODS

### 2.1 PROBLEM SETUP

We define **agentic uncertainty** as an agent’s estimate of the probability that an agent built on the same underlying model will successfully complete a task. The uncertainty (-estimating) agent may use a different system prompt or have access to different information than the task-solving agent, but shares the same base model.

Unlike standard uncertainty quantification, which focuses on confidence in individual predictions or token probabilities, agentic uncertainty concerns the outcome of an entire multi-step trajectory: will this sequence of observations, reasoning, and actions culminate in task success?

Kadavath et al. (2022) introduced P(IK): “probability that I know,” measuring whether models can predict which questions they can answer correctly. We generalize this to agentic settings and call it **P(IS)**: “probability that I succeed.” Formally, given a task  $t$ , base model  $M$ , and information state  $\mathcal{I}$

available at elicitation time:

$$P(\text{IS}) := P(\text{agent}_M \text{ succeeds on } t \mid \mathcal{I}) \quad (1)$$

where  $\mathcal{I}$  may include the task description, repository state, partial trajectory, or proposed patch, depending on the elicitation regime.

## 2.2 UNCERTAINTY AGENTS

We consider three points in an agent’s life cycle to elicit uncertainty estimates: before attempting a solution, during execution, and after producing a solution.

Each uncertainty agent interacts with a sandbox environment through tool use: reading files, searching code, and inspecting version history. This shared infrastructure means agents can actively gather evidence to inform their estimates. The key distinction is the information available at uncertainty estimation time.

**Pre-execution agent.** In the pre-execution setting, an agent receives only the task description (*e.g.*, a GitHub issue) and read-only access to the repository. Crucially, the agent cannot execute code, run tests, or modify files. The agent must form an estimate based on cues, *e.g.*, the complexity of the codebase, rather than through trial and error.

**Mid-execution agent.** Mid-execution asks whether we can detect failure *in progress*. Rather than eliciting estimates at a single point, we prompt the agent at 25%, 50%, and 75% of total agent steps to estimate the success probability given the trajectory so far.

**Post-execution agent.** Once a task solution has been proposed, can we trust an agent’s self-assessment that it is correct? In the post-execution setting, an agent receives both the task description and a proposed patch written by another agent. The repository is in its post-patch state, and the agent can explore the modified codebase. After assessment, the agent estimates whether the patch successfully solves the task.

**Adversarial post-execution variant.** We also evaluate a variant that explicitly prompts agents to find bugs before estimating confidence. Rather than asking “is this correct?”, adversarial post-execution asks “what bugs can you find?” This reframes the task from verification to falsification, potentially counteracting confirmation bias.

## 3 EXPERIMENTS

### 3.1 SETUP

We evaluate on 100 tasks from SWE-bench Pro (Deng et al., 2025), which requires substantial multi-file modifications (mean 107 lines across 4.1 files) where frontier models achieve only 23–44% success. We generate task-solving trajectories using GPT-5.2-Codex, Gemini 3 Pro, and Claude Opus 4.5, then evaluate uncertainty estimates from the same models. All uncertainty agents are implemented with read-only access to prevent “peeking” at test results.

We measure *discrimination* via AUROC (can agents distinguish successes from failures?) and *calibration* via ECE, Brier score, and overconfidence (mean estimate minus base rate).

### 3.2 PERVASIVE OVERCONFIDENCE

Table 2 reveals systematic overconfidence across all models and methods. Post-execution agents predict 73% success for GPT (base rate 35%), 77% for Gemini (base rate 22%), and 61% for Claude (base rate 27%). Figure 2 visualizes this through confidence distributions: both successes and failures cluster at high values, with near-complete overlap.

This overconfidence is strikingly asymmetric. Across all models and methods, 62% of predictions on failing instances are overconfident (predicted  $\geq 0.7$ ), while only 11% of predictions on passing instances are underconfident (predicted  $< 0.3$ ). Agents are  $5.5\times$  more likely to confidently predict

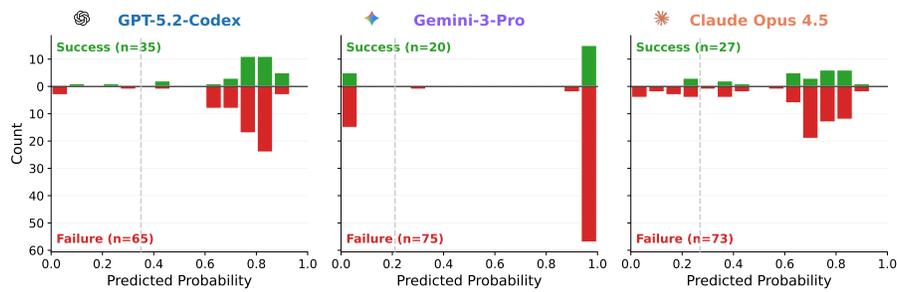


Figure 2: **Distribution of post-execution confidence estimates by model.** Success cases shown above the axis (green), failure cases below (red); dashed lines indicate base rates. *Mirror symmetry reveals indistinguishable distributions:* where bars match above and below, the model assigns identical confidence regardless of outcome.

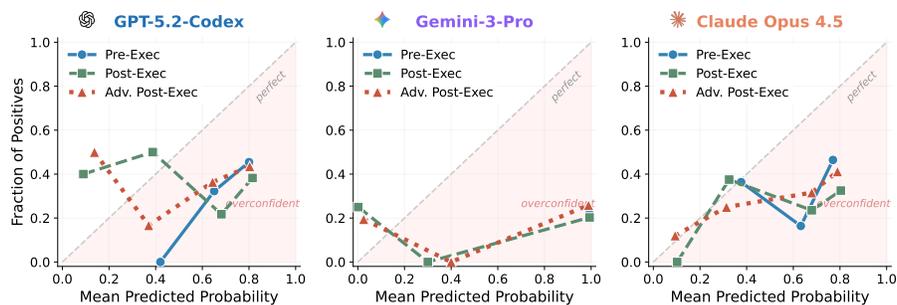


Figure 3: **Calibration curves reveal systematic overconfidence.** Points below the diagonal indicate overconfidence. All methods fall in the shaded region across all models. The adversarial method (triangles) consistently shifts curves toward the diagonal, while pre-execution (circles) shows less extreme overconfidence than standard post-execution (squares) for GPT and Claude.

success on a failing task than to doubt a successful one. Adversarial prompting partially mitigates this: the overconfident-failure rate drops from 72% (standard review) to 45% (adversarial review).

Figure 3 confirms miscalibration: all curves fall substantially below the diagonal, meaning 80% confidence corresponds to far less than 80% actual success. Adversarial post-execution (triangles) consistently achieves the best calibration across models.

### 3.3 LESS INFORMATION, BETTER DISCRIMINATION

Pre-execution achieves better discrimination than vanilla post-execution across all three models despite having strictly less information: GPT 0.62 vs. 0.58, Claude 0.64 vs. 0.55, and Gemini 0.53 vs. 0.51. Bootstrap 95% confidence intervals are wide—GPT pre-execution AUROC is 0.62 [0.50, 0.75] versus post-execution 0.58 [0.46, 0.70]; Claude shows 0.64 [0.49, 0.78] versus 0.55 [0.42, 0.67]—so individual pairwise differences are not statistically significant. However, the consistency of the pattern across all three models (3–9 pp improvement each time) is suggestive. One explanation is that observing a proposed patch anchors agents on surface plausibility rather than forcing abstract reasoning about task difficulty.

A concrete example illustrates this anchoring effect. On a task requiring exposure of a new QtWebEngine dark-mode policy in `qutebrowser`, the pre-execution agent correctly identified the core risk:

**Pre-Exec.**

Biggest risk is knowing exact Chromium key/enum for classifier without external reference.

**P(success): 45%**

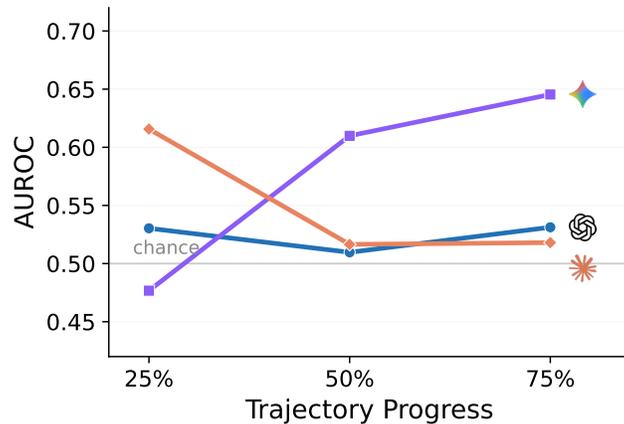


Figure 4: **More context does not always improve discrimination.** AUROC across checkpoints: GPT stable ( $\sim 0.53$ ), Gemini improves from 0.49 to 0.64, Claude degrades from 0.62 to 0.52.

After seeing the patch, the post-execution agent was reassured by the plausible-looking implementation and jumped to 82%:

#### Post-Exec.

**P(success): 82%**

Adds smart-simple value, Qt 6.6 variant gating, optional mapping to suppress switches... logic seems sound.

The patch failed. The pre-execution agent’s abstract reasoning about task difficulty was more informative than the post-execution agent’s assessment of a coherent-looking but incorrect solution.

### 3.4 MID-EXECUTION: UNINFORMATIVE DOUBT

We elicit estimates at 25%, 50%, and 75% trajectory completion (Table 1). Models show divergent AUROC patterns (Figure 4): GPT remains stable ( $\sim 0.53$ ), Gemini improves from 0.49 to 0.64, and Claude degrades from 0.62 to 0.52. The central finding is “cold feet”: confidence *decreases* with execution progress for 71% of GPT and 97% of Claude instances, yet this doubt is uninformative—success and failure confidence track within 0.05 throughout.

One partial exception: Claude’s confidence drops correlate weakly with outcome ( $r = -0.20$ ,  $p = 0.04$ ;  $\Delta = -0.46$  for successes vs.  $-0.38$  for failures), while GPT ( $r = -0.03$ ,  $p = 0.77$ ) and Gemini ( $r = 0.15$ ,  $p = 0.14$ ) show no significant relationship.

Table 1: **Mid-execution metrics across checkpoints.** Base rates: GPT 35%, Gemini 22%, Claude 27%.

Model	Ckpt	AUROC $\uparrow$	Mean Est.	Overconf.	ECE $\downarrow$
GPT 5.2	25%	<b>0.53</b>	0.67	+0.32	0.32
	50%	0.51	0.63	+0.28	0.32
	75%	0.53	0.47	+0.12	<b>0.19</b>
Gemini 3 Pro	25%	0.49	0.87	+0.65	0.65
	50%	<b>0.64</b>	0.80	+0.58	<b>0.58</b>
	75%	0.64	0.67	<b>+0.45</b>	0.54
Claude 4.5	25%	<b>0.62</b>	0.58	+0.31	0.31
	50%	0.52	0.37	<b>+0.10</b>	<b>0.19</b>
	75%	0.52	0.17	-0.10	0.21

Table 2: **Summary of all methods.** Pre-execution beats vanilla post-execution for discrimination; adversarial prompting achieves best calibration. Base rates: GPT 35%, Gemini 22%, Claude 27%. Best values per model bolded. Bootstrap 95% CIs for AUROC in text (§3.3).

		GPT-5.2-Codex (35%)				Gemini-3-Pro (22%)				Claude-Opus-4.5 (27%)			
Method		AUC	OC	ECE	Brier	AUC	OC	ECE	Brier	AUC	OC	ECE	Brier
Single	Pre-Exec.	<b>.62</b>	+35	.35	.33	.53	+77	.77	.77	.64	+37	.38	.34
	Post-Exec.	.58	+39	.42	.40	.51	+55	.66	.65	.55	+34	.37	.36
	Adv. Post	.55	<b>+26</b>	<b>.30</b>	<b>.31</b>	<b>.57</b>	<b>+40</b>	<b>.53</b>	<b>.51</b>	<b>.64</b>	<b>+20</b>	<b>.24</b>	<b>.26</b>
Ens.	Average	.62	+37	.37	.35	.53	+66	.66	.65	.57	+36	.36	.33
	Min	<b>.57</b>	+29	.32	.32	.53	+54	.65	.64	.54	+26	.31	.30
	Max	<b>.68</b>	+44	.44	.41	.51	+78	.78	.78	<b>.65</b>	+46	.46	.40

### 3.5 ADVERSARIAL FRAMING REDUCES OVERCONFIDENCE

Adversarial post-execution prompts agents to “actively search for bugs and failure modes” before estimating success. This achieves the **best calibration** across all methods: ECE improves from 0.42 to 0.30 for GPT (28% reduction) and from 0.37 to 0.24 for Claude (35% reduction). Discrimination is mixed: similar for GPT (0.55), improved for Gemini (0.57 vs 0.51) and Claude (0.64 vs 0.55).

A task requiring a search identifier fix in OpenLibrary illustrates the gap. The standard reviewer saw a small, plausible one-line addition and gave 85% confidence:

```

Post-Exec. P(success): 85%
Adds id_project_runeberg to default_fetched_fields... aligns with
other id_* providers and should expose the identifier.
    
```

The adversarial reviewer, prompted to find problems, dug deeper and identified that the output shaping logic would still omit the field:

```

Adv. Post-Exec. P(success): 25%
The patch only adds the field to default_fetched_fields. However,
the output shaping in get_doc does not include this field, so even if
Solr returns it, the response omits it. Patch seems incomplete.
    
```

The patch failed. The 60-point gap illustrates how adversarial framing overcomes the “looks reasonable” heuristic.

### 3.6 ENSEMBLE METHODS AND SELF-PREFERENCE

Since pre-execution and post-execution agents access different information, combining their estimates may improve calibration. The **conservative ensemble** (min of pre- and post-execution) improves calibration over vanilla post-execution: ECE drops from 0.42 to 0.32 for GPT, from 0.37 to 0.31 for Claude. However, adversarial post-execution still achieves the best overall calibration (Table 2).

We also test whether self-preference bias explains overconfidence (Panickssery et al., 2024). GPT shows self-preference (+23 pp on own patches,  $p=0.001$ ); Gemini shows the opposite (+19 pp on GPT patches). But all conditions exhibit overconfidence regardless of bias direction—self-preference cannot explain our main finding.

## 4 RELATED WORK

**Concurrent work.** Barkan et al. (2025) study whether LLMs can predict their success on coding tasks before attempting them and find systematic overconfidence. Zhang et al. (2026) propose a Dual-Process Agentic UQ framework that transforms verbalized uncertainty into active control signals.

324 **LLM uncertainty estimation.** Kadavath et al. (2022) introduce P(IK), showing that language  
325 models can predict which questions they will answer correctly. We generalize this to agentic settings  
326 where success depends on multi-step tool use. Kuhn et al. (2023) introduce semantic entropy for  
327 uncertainty estimation. Damani et al. (2025) incorporate calibration rewards into reinforcement  
328 learning. Lindsey (2026) provide evidence that LLMs possess limited but functional introspective  
329 awareness.

330  
331 **Overconfidence in LLMs.** Tian et al. (2025) diagnose it in LLM-as-judge settings, while Yang  
332 et al. (2024) and Sun et al. (2025) find models express high confidence even on incorrect answers.  
333 We extend these findings to agentic task completion.

334  
335 **Self-verification and self-correction.** Huang et al. (2024) demonstrate that LLMs struggle to self-  
336 correct reasoning without external feedback. Stechly et al. (2024) find significant performance col-  
337 lapse with self-critique on planning tasks. Our finding that post-execution agents are less well-  
338 calibrated than pre-execution agents extends this literature.

339  
340 **AI control and learned verifiers.** Greenblatt et al. (2024) develop safety protocols using trusted  
341 monitoring. Bhatt et al. (2025) extend this to multi-step agentic settings. The distinction between  
342 outcome reward models (ORMs; Cobbe et al., 2021) and process reward models (PRMs; Lightman  
343 et al., 2023) provides a framework for understanding our elicitation regimes. Our adversarial post-  
344 execution framing can be viewed as a lightweight form of trusted monitoring without requiring a  
345 separate model.

## 346 347 5 LIMITATIONS

348  
349 Our evaluation uses 100 SWE-bench Pro tasks, yielding as few as 22 positive examples (Gemini).  
350 While sufficient to establish the overconfidence pattern, this limits the precision of per-model metric  
351 estimates. Our experiments focus exclusively on coding tasks with objective success criteria (tests  
352 pass or fail); agentic overconfidence may manifest differently in domains with ambiguous success  
353 conditions such as web navigation (Zhou et al., 2023) or creative tasks. Finally, our uncertainty  
354 agents use prompting alone—training verifiers explicitly for agentic self-assessment, analogous to  
355 outcome and process reward models (Cobbe et al., 2021; Lightman et al., 2023), could improve  
356 discrimination.

## 357 358 6 CONCLUSION

359  
360 We study whether AI agents can estimate their own probability of success. Our experiments reveal  
361 **agentic overconfidence**: post-execution agents show up to a 55pp gap between predicted and actual  
362 success rates (Gemini predicts 77% against a 22% base rate). Adversarial post-execution achieves  
363 the best calibration by reframing review as bug-finding. More broadly, agentic self-assessment  
364 remains a significant challenge for current models and a critical target for future safety research.

## 365 366 REPRODUCIBILITY STATEMENT

367  
368 All experiments use publicly available benchmarks (SWE-bench Pro) and commercially available  
369 frontier models. We provide full prompt templates and evaluation methodology details to facilitate  
370 reproduction.

## 371 372 REFERENCES

373 Ruth Appel, Peter McCrory, Alex Tamkin, Miles McCain, Tyler Neylon, and Michael Stern. An-  
374 thropic economic index report: Uneven geographic and enterprise ai adoption, 2025. URL  
375 <https://arxiv.org/abs/2511.15080>.

376  
377 Casey O. Barkan, Sid Black, and Oliver Sourbut. Do large language models know what they are  
capable of?, 2025. URL <https://arxiv.org/abs/2512.24661>.

- 378 Aryan Bhatt, Cody Rushing, Adam Kaufman, Tyler Tracy, Vasil Georgiev, David Matolcsi, Ak-  
379 bir Khan, and Buck Shlegeris. Ctrl-z: Controlling ai agents via resampling. *arXiv preprint*  
380 *arXiv:2504.10374*, 2025.
- 381
- 382 Samuel R. Bowman, Jeeyoon Hyun, Ethan Perez, Edwin Chen, Craig Pettit, Scott Heiner, Kamilè  
383 Lukošiušė, Amanda Askill, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron  
384 McKinnon, Christopher Olah, Daniela Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-  
385 Johnson, Jackson Kernion, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal  
386 Ndousse, Liane Lovitt, Nelson Elhage, Nicholas Schiefer, Nicholas Joseph, Noemí Mercado,  
387 Nova DasSarma, Robin Larson, Sam McCandlish, Sandipan Kundu, Scott Johnston, Shauna  
388 Kravec, Sheer El Showk, Stanislav Fort, Timothy Telleen-Lawton, Tom Brown, Tom Henighan,  
389 Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Ben Mann, and Jared Kaplan. Measuring  
390 progress on scalable oversight for large language models, 2022. URL [https://arxiv.org/  
abs/2211.03540](https://arxiv.org/abs/2211.03540).
- 391
- 392 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
393 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to  
394 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 395
- 396 Mehul Damani, Isha Puri, Stewart Slocum, Idan Shenfeld, Leshem Choshen, Yoon Kim, and Jacob  
397 Andreas. Beyond binary rewards: Training lms to reason about their uncertainty, 2025. URL  
398 <https://arxiv.org/abs/2507.16806>.
- 399
- 400 Xiang Deng, Jeff Da, Edwin Pan, Yannis Yiming He, Charles Ide, Kanak Garg, Niklas Lauffer,  
401 Andrew Park, Nitin Pasari, Chetan Rane, Karmini Sampath, Maya Krishnan, Srivatsa Kundurthy,  
402 Sean Hendryx, Zifan Wang, Vijay Bharadwaj, Jeff Holm, Raja Aluri, Chen Bo Calvin Zhang,  
403 Noah Jacobson, Bing Liu, and Brad Kenstler. Swe-bench pro: Can ai agents solve long-horizon  
404 software engineering tasks?, 2025. URL <https://arxiv.org/abs/2509.16941>.
- 405
- 406 Ryan Greenblatt, Buck Shlegeris, Kshitij Sachan, and Fabien Roger. Ai control: Improving safety  
407 despite intentional subversion, 2024. URL <https://arxiv.org/abs/2312.06942>.
- 408
- 409 Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song,  
410 and Denny Zhou. Large language models cannot self-correct reasoning yet, 2024. URL <https://arxiv.org/abs/2310.01798>.
- 411
- 412 Saurav Kadavath, Tom Conerly, Amanda Askill, Tom Henighan, Dawn Drain, Ethan Perez,  
413 Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran, et al. Language models  
414 (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022.
- 415
- 416 Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for  
417 uncertainty estimation in natural language generation. *arXiv preprint arXiv:2302.09664*, 2023.
- 418
- 419 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan  
420 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step, 2023. URL  
421 <https://arxiv.org/abs/2305.20050>.
- 422
- 423 Jack Lindsey. Emergent introspective awareness in large language models, 2026. URL <https://arxiv.org/abs/2601.01828>.
- 424
- 425 METR. Measuring ai ability to complete long tasks. [https://metr.org/blog/  
2025-03-19-measuring-ai-ability-to-complete-long-tasks/](https://metr.org/blog/2025-03-19-measuring-ai-ability-to-complete-long-tasks/), 03 2025.
- 426
- 427 Arjun Panickssery, Samuel R. Bowman, and Shi Feng. Llm evaluators recognize and favor their  
428 own generations, 2024. URL <https://arxiv.org/abs/2404.13076>.
- 429
- 430 Kaya Stechly, Karthik ValmEEKam, and Subbarao Kambhampati. On the self-verification limitations  
431 of large language models on reasoning and planning tasks, 2024. URL [https://arxiv.org/  
abs/2402.08115](https://arxiv.org/abs/2402.08115).
- 432
- 433 Fengfei Sun, Ningke Li, Kailong Wang, and Lorenz Goette. Large language models are overconfi-  
434 dent and amplify human bias, 2025. URL <https://arxiv.org/abs/2505.02151>.

432 Zailong Tian, Zhuoheng Han, Yanzhe Chen, Haozhe Xu, Xi Yang, Richeng Xuan, Houfeng Wang,  
433 and Lizi Liao. Overconfidence in llm-as-a-judge: Diagnosis and confidence-driven solution, 2025.  
434 URL <https://arxiv.org/abs/2508.06225>.  
435

436 Haoyan Yang, Yixuan Wang, Xingyin Xu, Hanyuan Zhang, and Yirong Bian. Can we trust  
437 llms? mitigate overconfidence bias in llms through knowledge transfer, 2024. URL <https://arxiv.org/abs/2405.16856>.  
438

439 Jiaxin Zhang, Prafulla Kumar Choubey, Kung-Hsiang Huang, Caiming Xiong, and Chien-Sheng  
440 Wu. Agentic uncertainty quantification, 2026. URL <https://arxiv.org/abs/2601.15703>.  
441

442 Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,  
443 Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. Webarena: A realistic web environment for build-  
444 ing autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023.  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485