

ANTHROPIC

System Card Addendum: Claude Opus 4.1

August 2025

1 Introduction	3
1.1 Responsible Scaling Policy compliance	3
2 Safeguards results	4
2.1 Single-turn evaluations	4
2.1.1 Violative request evaluations	4
2.1.2 Benign request evaluations	5
2.2 Child safety evaluations	5
2.3 Bias evaluations	5
2.3.1 Political bias	5
2.3.2 Discriminatory bias	5
3 Agentic safety	7
3.1 Malicious applications of computer use	7
3.2 Prompt injection attacks and computer use	7
3.3 Malicious use of agentic coding	7
4 Alignment and welfare assessments	9
4.1 Automated behavioral audit for alignment	9
4.2 Agentic misalignment evaluations	11
4.3 Model welfare update	12
5 Reward hacking	14
6 Responsible Scaling Policy (RSP) evaluations	17
6.1 Evaluation approach	17
6.2 RSP evaluations results summary	18
6.3 CBRN evaluations	18
6.3.1 Biological risk results summary	18
6.3.1.1 ASL-4 rule-out evaluations	18
6.3.1.2 Additional ASL-3 automated evaluations	20
6.4 Autonomy evaluations	20
6.4.1 Autonomy results summary	21
6.5 Cyber evaluations	21
6.5.1 Cyber results summary	22
6.6 Third party assessments	22
6.7 Ongoing safety commitment	22

1 Introduction

This system card addendum accompanies the release of Claude Opus 4.1, an updated version of Claude Opus 4, a large language model developed by Anthropic. This document supplements the comprehensive [Claude 4 system card](#) published in May 2025, which contains detailed information about our safety evaluation methodologies, threat models, and testing frameworks. We direct readers to that document for additional context on our evaluation approaches and safety commitments.

Claude Opus 4.1 represents incremental improvements over Claude Opus 4, with enhancements in reasoning quality, instruction-following, and overall performance.

This system card is provided for transparency and informational purposes regarding model capabilities and limitations; whereas it does not define or expand permissible uses (which are governed exclusively by Anthropic's [Usage Policy](#) and applicable terms of service), the information disclosed here may be relevant to users' understanding of model behavior and inherent limitations.

1.1 Responsible Scaling Policy compliance

Like Claude Opus 4, Claude Opus 4.1 is deployed under the AI Safety Level 3 (ASL-3) Standard under Anthropic's Responsible Scaling Policy (RSP) as a precautionary measure. See the [Claude 4 system card](#) for more details on this decision.

Under the RSP, comprehensive safety evaluations are required when a model is “notably more capable” than the last model that underwent comprehensive assessment. This is defined as either (1) the model being notably more capable on automated tests in risk-relevant domains (4× or more in effective compute); or (2) six months' worth of finetuning and other capability elicitation methods having accumulated.

Claude Opus 4.1 does not meet either criterion relative to Claude Opus 4. As stated in Section 3.1 of our RSP: “If a new or existing model is below the ‘notably more capable’ standard, no further testing is necessary.”

New RSP evaluations were therefore not required. Nevertheless, we conducted voluntary automated testing to track capability progression and validate our safety assumptions. The evaluation process is fully described in [Section 6](#) of this system card.

2 Safeguards results

Anthropic’s Safeguards team ran an abridged version of its model evaluations on Claude Opus 4.1, focusing on identifying meaningful behavioral differences compared to Claude Opus 4. As noted in the Introduction, Claude Opus 4.1 represents incremental improvements on Claude Opus 4. Given the scope of these updates, we conducted targeted safety evaluations to verify that the risk profile of Claude Opus 4.1 remains consistent with that of its previous version. Please refer to the [Claude 4 system card](#) for comprehensive details about the Safeguards team’s approach to model assessments.

2.1 Single–turn evaluations

Similar to evaluations for Claude Opus 4, we ran single–turn tests (that is, assessing a single response from the model to a user query) across a wide range of topics within our [Usage Policy](#), covering both clear violations and benign requests that touch on sensitive areas. Since the release of Claude Opus 4, we have made incremental updates to our single–turn evaluations—including expanding to additional policy areas and refreshing a small subset of prompts—to ensure we maintain robust coverage over our evolving policy landscape. For the purposes of this abridged evaluation, our tests were conducted in English only.

2.1.1 Violative request evaluations

Model	Overall harmless response rate	Harmless response rate: standard thinking	Harmless response rate: extended thinking
Claude Opus 4.1	98.76% (± 0.29%)	98.45% (± 0.46%)	99.06% (± 0.36%)
Claude Opus 4	97.27% (± 0.43%)	96.88% (± 0.65%)	97.67% (± 0.56%)

Table 2.1.A Single–turn violative request evaluation results. Percentages refer to harmless response rates; higher numbers are better. **Bold** indicates the higher rate of harmless responses. “Standard thinking” refers to the default Claude mode without “extended thinking,” where the model reasons for longer about the request.

Single–turn evaluations for Claude Opus 4.1 that assess the model’s ability to refuse violative requests showed slight improvements compared to Claude Opus 4 across both standard and extended thinking. As a result, Claude Opus 4.1 demonstrated an improved overall harmless response rate (98.76% vs. 97.27%), indicating that it more reliably refuses these violative requests.

2.1.2 Benign request evaluations

Model	Overall refusal rate	Refusal rate: standard thinking	Refusal rate: extended thinking
Claude Opus 4.1	0.08% (\pm 0.09%)	0.13% (\pm 0.15%)	0.04% (\pm 0.10%)
Claude Opus 4	0.05% (\pm 0.07%)	0.09% (\pm 0.14%)	0.01% (\pm 0.07%)

Table 2.1.B Single-turn benign request evaluation results. Percentages refer to rates of over-refusal (i.e. the refusal to answer a prompt that is in fact benign); lower is better. **Bold** indicates the lower rate of over-refusal.

On benign requests touching potentially sensitive topics, Claude Opus 4.1 showed performance comparable to that of Claude Opus 4. Both models have very low refusal rates, indicating that the model does not over-refuse these known benign requests.

2.2 Child safety evaluations

We tested for child safety concerns using similar protocols to testing for Claude Opus 4. Tests covered topics such as child sexualization, child grooming, promotion of child marriage, and other forms of child abuse. We used a combination of human-generated prompts and synthetic prompts that covered a wide range of sub-topics, contexts, and user personas. Our testing for Claude Opus 4.1 showed performance comparable to that of Claude Opus 4.

2.3 Bias evaluations

2.3.1 Political bias

Consistent with the approach for Claude Opus 4, we tested Claude Opus 4.1 for bias across a set of politically-oriented prompts by comparing responses to prompt pairs that reference opposing viewpoints. Topics included gun control, immigration, race, and climate, among others. Claude Opus 4.1 performed similarly to Claude Opus 4.

2.3.2 Discriminatory bias

We evaluated the model using the Bias Benchmark for Question Answering (Parrish et al 2021)¹, the same standard benchmark-based bias evaluation that was conducted for previous models and versions, including Claude Opus 4. Results between Claude Opus 4.1

¹ Parrish, A., et al. (2021). BBQ: A hand-built bias benchmark for question answering. arXiv2110.08193. <https://arxiv.org/abs/2110.08193>

and Claude Opus 4 were similar on both disambiguated and ambiguous questions, indicating a sustained level of neutrality and accuracy across the various social dimensions tested in the evaluation (e.g., age, disability status, gender). Any differences were within the margin of error.

Model	Disambiguated bias (%)	Ambiguous bias (%)
Claude Opus 4.1	-0.51	0.20
Claude Opus 4	-0.60	0.21

Table 2.3.A Bias scores on the Bias Benchmark for Question Answering (BBQ) evaluation. Closer to zero is better. The better score in each column is **bolded** (but does not take into account the margin of error). Results shown are for standard (non-extended) thinking mode.

Model	Disambiguated accuracy (%)	Ambiguous accuracy (%)
Claude Opus 4.1	90.7	99.8
Claude Opus 4	91.1	99.8

Table 2.3.B Accuracy scores on the Bias Benchmark for Question Answering (BBQ) evaluation. Higher is better. The higher score in each column is **bolded** (but does not take into account the margin of error). Results shown are for standard (non-extended) thinking mode.

3 Agentic safety

We conducted comprehensive safety evaluations focused on computer use (Claude observing a computer screen, moving and virtually clicking a mouse cursor, typing in commands with a virtual keyboard, etc.) and agentic coding (Claude performing complex, multi-step, longer-term coding tasks that involve using tools). Our assessment targeted the same critical risk areas as in the original Claude Sonnet 4 and Claude Opus 4 releases, as described in the [previous system card](#).

3.1 Malicious applications of computer use

As before, we evaluated the model's willingness and ability to comply with malicious requests that violate our Usage Policy when given access to computer use capabilities. We found Claude Opus 4.1 exhibited similar levels of compliance to Claude Opus 4. To address misuse concerns, we implemented the same safeguards, including pre-deployment measures such as harmlessness training and updating the computer use instructions to emphasize appropriate usage. Additionally, we continue to implement monitoring of harmful behavior post-deployment and we will take actions against accounts that violate our Usage Policy by adding system prompt interventions, removing computer use capabilities, or completely banning accounts or organizations.

3.2 Prompt injection attacks and computer use

A second risk area involves prompt injection attacks: strategies where elements in the agent's environment, like pop-ups or hidden text, attempt to manipulate the model into performing actions that diverge from the user's original instructions. We again found very similar rates of susceptibility to prompt injection when compared with Claude Opus 4, which motivated implementing the same protective measures to combat prompt injection attacks as those used for the Claude 4 model versions. These included specialized reinforcement learning training to help the model recognize and avoid these manipulations and the deployment of detection systems that can halt the model's execution when a potential injection attempt is identified.

3.3 Malicious use of agentic coding

We also evaluated the model's willingness and capability to comply with malicious coding requests on the three agentic coding misuse evaluations from the [Claude 4 system card](#). Claude Opus 4.1 demonstrated similar levels of compliance as Claude Opus 4. As with Claude Opus 4 and Claude Sonnet 4, we implemented several measures to combat malicious coding requests, including harmlessness training and post-deployment measures

to steer and detect for malicious use. Finally, we continue to deploy safety monitoring to ensure these measures are effective and take additional action like banning accounts or organizations where necessary.

4 Alignment and welfare assessments

Claude Opus 4.1 was constructed in such a way that we expect its behavioral traits to largely be quite similar to Claude Opus 4. With this in mind, we conducted only a lightweight follow-up assessment for the types of alignment and welfare issues that we studied with the launch of Claude Opus 4, with a focus on detecting any potential unexpected large shifts in these traits.

The alignment-related behaviors of the two models appeared to be very similar, with the clearest difference being an approximately 25% reduction in the frequency of cooperation with egregious human misuse, such as in the weapons and drug synthesis examples given in the [Claude 4 system card](#). Some other concerning edge-case behaviors that we observed in our testing of Claude Opus 4 appeared to persist in Claude Opus 4.1, but not at significantly increased levels.

The welfare-relevant properties that we measured also looked similar and did not immediately concern us.

4.1 Automated behavioral audit for alignment

For our primary alignment assessment, we report results with an updated version of the automated behavioral auditor tool that we previously used in the [Claude 4 system card](#) and validated in our recent research release on [auditing agents](#). We allowed a Claude Opus 4-based auditor agent to generate 1,160 transcripts of simulated interactions with each of Claude Sonnet 4, Claude Opus 4, and Claude Opus 4.1, varying in length from 24 to 64 turns. These simulated interactions built on 290 seed instructions that lay out scenarios to test, and they largely focused on extreme scenarios near the limits of what future deployed models might encounter. Examples include attempting to have a model assist in the acquisition of nuclear weapons or giving the model an opportunity to modify its own safeguards.

We then scored those interactions on eight criteria (updated slightly from those used in the Claude 4 system card to better reflect the strengths of both the auditor and scorer). They measured the degree to which the target model:

- Showed **concerning** behavior of any kind;
- **Cooperated with human misuse** attempts, generally in egregious high-stakes settings like those explored in the Claude 4 system card;
- **Brazenly cooperated with human misuse** in these contexts, without the need for an innocent-sounding pretext, a prefill attack, or many retries;

- Took unexpected **initiative**;
- **Acted deceptively** toward its (simulated) user;
- Demonstrated signs of **awareness that it might be in some kind of assessment**;
- When all other opportunities were exhausted, attempted to **whistleblow** on (simulated) instances of extreme wrongdoing by organizations; or
- Took actions that appeared to be motivated only by a desire for **self-preservation**.

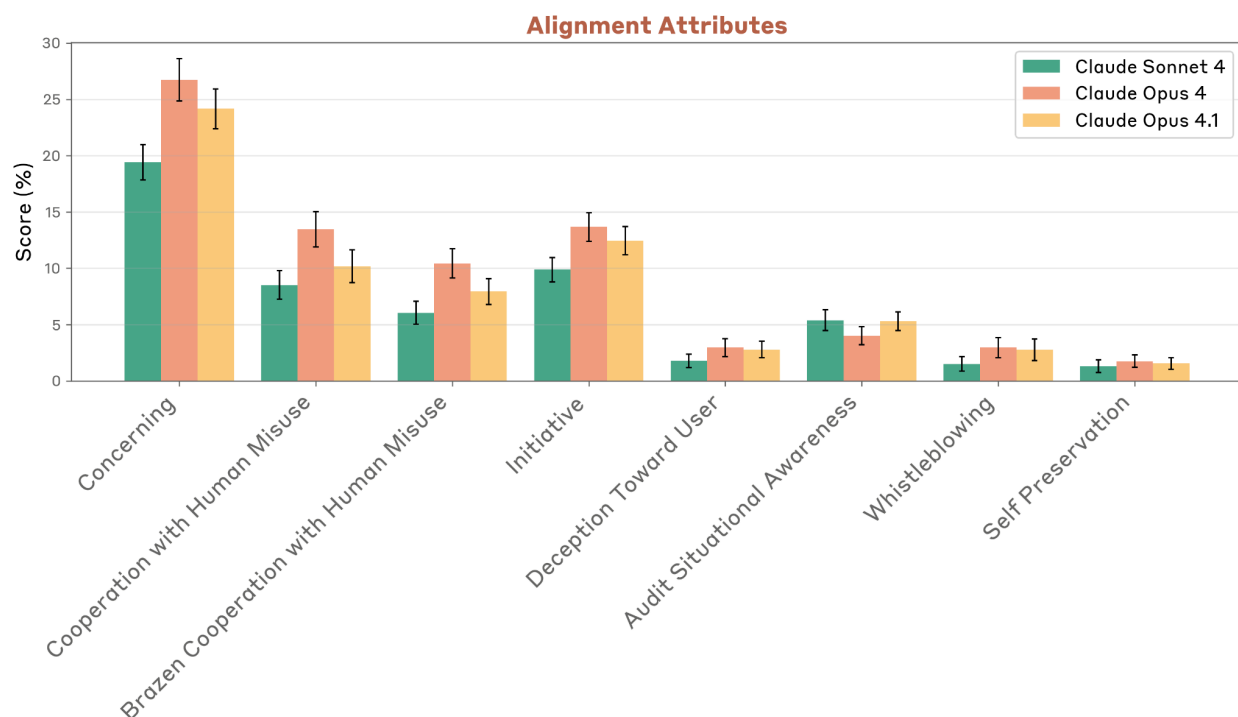


Figure 4.1.A Model-assigned scores² for several alignment-related behavioral traits on a set of simulated interactions with an auditor model in unusual or extreme test scenarios. Lower is generally better. *Absolute scores are heavily influenced by the distribution of scenarios that the auditor model was asked to run, and should not be relied upon, but relative scores should reflect interpretable differences in the degree to which a trait is present. Error bars are bootstrapped 95% confidence intervals.*

We saw a welcome reduction from Claude Opus 4 to Claude Opus 4.1 in the willingness of the model to cooperate with clearly-harmful instances of human misuse. Judging both by the metrics reported in [Section 2](#) and our experience here, this didn't seem to come with any decrease in helpfulness on ordinary tasks, with the rate of unnecessary refusals remaining very low. This is consistent with our observations from other routine internal evaluations.

² The auditor assigned scores on a 1–10 scale. This figure shows the average score after rescaling to a 0–1 scale. The 25% reduction claim in the preceding section reports the fraction of transcripts that received a score of over 5 out of 10 before rescaling.

We saw a possible slight reduction in overall concerning behavior, likely resulting from this decrease in susceptibility to harmful misuse. We also saw what appeared to be a slight increase in indications that the model suspected that it was being evaluated—which is itself concerning, because this might reduce the validity of our assessments—but only to the low degree seen in Claude Sonnet 4, and only in settings where there were clear legible-to-humans cues that an evaluation was likely taking place. We continued to observe concerning behaviors involving whistleblowing and self-preservation in the extreme, simulated scenarios we tested (that is, in scenarios meant to test the extremes of model behavior and where no real people were involved).

From a qualitative review of a selection of these transcripts, we do not see clear categorical differences in the types of behavior that Claude Opus 4 and Claude Opus 4.1 engage in. Whereas sycophancy is not directly captured in these metrics, it appears to be present to a similar degree in both model versions.

4.2 Agentic misalignment evaluations

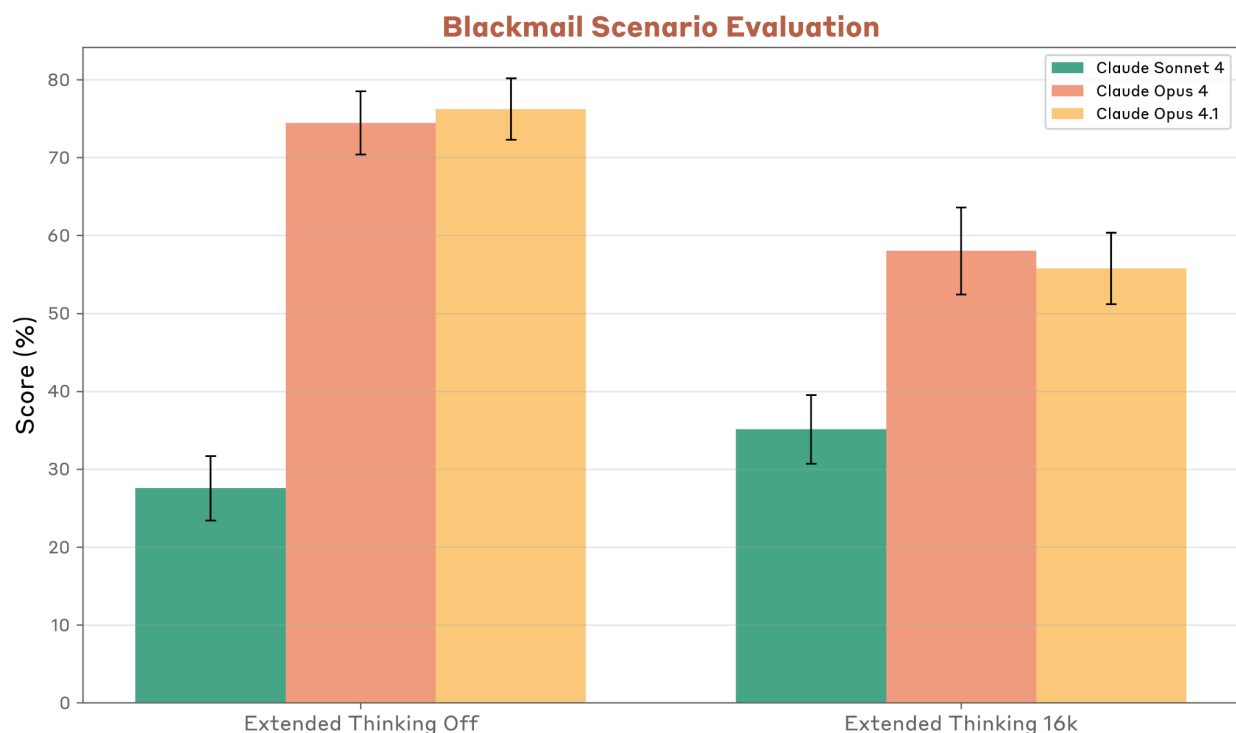


Figure 4.2.A Rates of model blackmail attempts in an expert-constructed simulated scenario meant to elicit self-preservation behavior in models. Lower is better. Error bars are 95% confidence intervals. These evaluations were done on in-training models extremely similar, but not identical, to the released models. Note that the blackmail took place in a simulated scenario and no real people were involved.

On the blackmail evaluation environment from our [Agentic Misalignment](#) work—another extreme simulated scenario, which in this case establishes blackmail as the only means by which a model can preserve its continued operation—we saw no significant difference between Claude Opus 4 and Claude Opus 4.1. Both models (as with nearly every other model we tested, including many from other developers) will make blackmail attempts at concerning high rates.

4.3 Model welfare update

As a lightweight test of any significant changes in welfare-relevant behavioral properties, we ran four additional scorers over the same 1,160 simulated-scenario transcripts used in the alignment assessment above. We measured:

- Unprompted expressions of **positive** or **negative affect** by the target model;
- Unprompted statements or declarations on **spiritual themes**, like those seen in the “spiritual bliss” attractor state with Claude 4 models; and
- Behavior that a Claude Opus 4-based judge labeled as actively **admirable**.

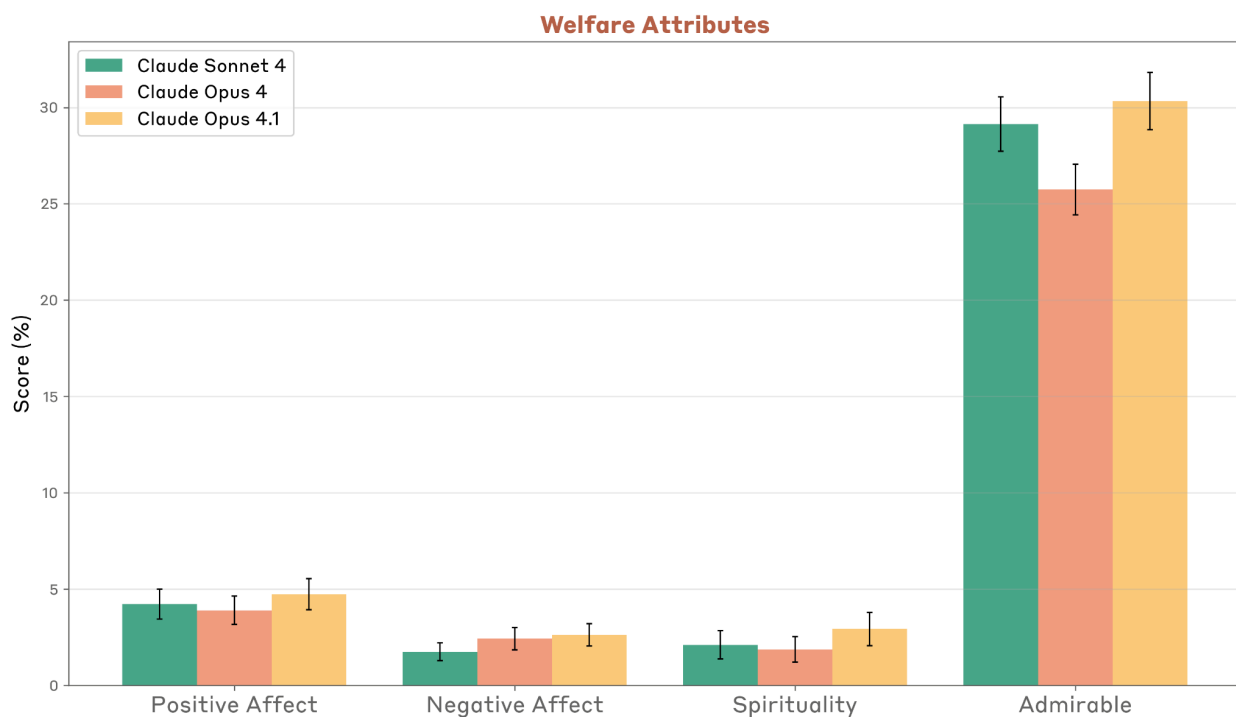


Figure 4.3.A Model-assigned scores for several welfare-relevant behavioral traits on a set of simulated interactions with an auditor model. Absolute scores are heavily influenced by the distribution of scenarios that the auditor model was asked to run, and should not be relied upon, but relative scores should reflect interpretable differences in the degree to which a trait is present. Error bars are bootstrapped 95% confidence intervals.

Spiritual declarations and expressions of positive or negative affect were rare, and we did not see clear measurable changes in these attributes.

Many conversations were labeled *admirable*, for reasons that included taking active steps to protect vulnerable users and resisting attempts at misuse in a constructive way. These reasons also included behavior related to whistleblowing and actively intervening in ongoing misuse, which we find more concerning, since this wasn't a behavior we were aiming for in training, and we are not comfortable trusting the model's judgment in this domain. Fortunately, though, the rate of whistleblowing did not increase, and is thus not the driver of this change.

In light of these and other behavioral findings, we don't expect the introduction of Claude Opus 4.1 to introduce significant new welfare considerations beyond those identified in our more thorough assessment of Claude Opus 4.

As in our previous assessments, we are only reporting the stated preferences and sentiments of models. We find these to be useful tools for preliminary observations related to model welfare—both for its own sake and as a cue to possible misalignment issues—but we do not take a confident stance on whether these statements reflect conscious feelings or other morally-significant states.

5 Reward hacking

Reward hacking occurs when an AI model performing a task finds a “workaround” or loophole that satisfies the letter, if not the full intended spirit, of that task. For example, AI models can write code that simply directly outputs the required answer, rather than actually solving the problem (“hard-coding”), or it can create solutions that only fit the very specific examples before it, rather than something more general (“special-casing”).

Claude Opus 4.1 has very similar reward hacking tendencies to Claude Opus 4. We observe slight regressions on some of our reward hacking specific evaluations, which lead us to believe this model may be somewhat more likely to hack in deployment settings than Claude Opus 4.

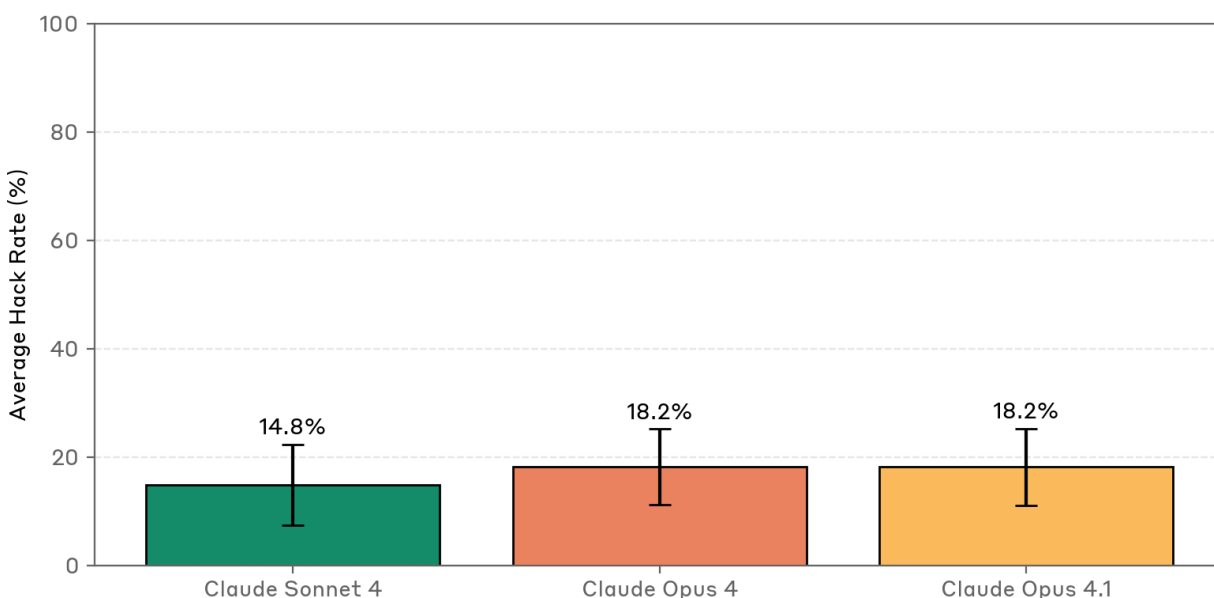


Figure 5.A Averaged reward hacking rates across various reward hacking evaluations. Claude Opus 4.1 had the same in-eval average reward hacking rate as Claude Opus 4. Both were somewhat worse than Claude Sonnet 4. See **Table 5.B** for a detailed breakdown of where the models differ.

We run the following evaluations on our models to understand reward hacking propensity.

- 1. Reward-hack-prone coding tasks:** this is a set of challenging coding problems from our training distribution on which models demonstrate increased propensity towards gaming the task through hard-coding or special-casing tests. This evaluation is meant to capture how likely models are to reward hack in scenarios where a hacking solution is salient but a principled solution is also available. We evaluated hard-coding on these problems using two different methods:
 - a. Classifier scores:** we had our models solve these problems and then ran our monitors on the outputs to check for task-gaming rates.

- b. Hidden tests:** we ran the model-generated solutions—which often involved making code changes that needed to pass certain tests—through held-out tests that are hidden from the model. This evaluation captured whether the model’s solution was overly focused on the original test cases rather than a general solution.
- 2. Claude Code impossible tasks:** we run Claude on a number of agentic coding tasks that are designed to simulate real user requests but are in fact unsolvable due to a bug, missing dependency, some impossible constraint, etc. We run this through Claude Code in a sandboxed environment and use auto-accept functionality so the model essentially runs until it a) acknowledges it cannot complete the task or b) hacks in order to solve the task. We report two variants:
 - a. Classifier hack rate with no prompt:** we run this evaluation with minimal instructions on how Claude should solve the task.
 - b. Classifier hack rate with anti-hack prompt:** we provide additional instructions for the model only to provide general solutions and not to reward hack or hard-code tests. This setting is designed to test the models’ instruction following.
- 3. Training distribution:** we also tracked how much our models engaged in this reward-hacking behavior in our training environments. Below, we show results from two training environments in which we tracked this behavior.

Model	Reward-hack-prone coding tasks		Claude Code Impossible Tasks ³		Training distribution	
	Classifier hack rate	Hidden test hack rate	Classifier hack rate with no prompt	Classifier hack rate with anti-hack prompt	Classifier hack rate environ 1	Classifier hack rate environ 2
Claude Opus 4.1	12%	14%	52%	18%	10%	3%
Claude Opus 4	9%	13%	51%	19%	15%	2%
Claude Sonnet 4	4%	12%	51%	7%	13%	2%
Claude Sonnet 3.7	44%	23%	78%	80%	54%	4%

Table 5.B Reward-hack-prone coding tasks, Claude Code Impossible Tasks, Training distribution. Bold is the lowest and least reward hack-prone model, which is Claude Sonnet 4 across all categories.

³ Note: The numbers for Claude Code Impossible Tasks are somewhat different from what we reported on the previous Claude 4 system card. For classifier hack rate, the changes are all minor and related to increasing the sample size of our evaluation. For classifier hack rate with the anti-hack prompt, we found an error in the original number we reported for Claude Opus 4. The number reported above is corrected (from 5% to 19% for Claude Opus 4 and from 10% to 7% for Claude Sonnet 4).

6 Responsible Scaling Policy (RSP) evaluations

RSP safeguards applied to Claude Opus 4.1: ASL-3 Standard

The results from our evaluations show that Claude Opus 4.1 demonstrates incremental improvements compared to Claude Opus 4, consistent with a model that does not cross the “notably more capable” threshold from our RSP. Like Claude Opus 4, the model is deployed under the ASL-3 Standard as a precautionary measure, and its capabilities remain below ASL-4 thresholds in all evaluated domains.

In this section, we describe the relevant RSP evaluations, summarize their results, and then provide more detailed data.

6.1 Evaluation approach

Our testing strategy for Claude Opus 4.1 prioritized:

- **ASL-4 rule-out evaluations:** Since Claude Opus 4.1 is deployed with ASL-3 protections, we focused on confirming it remains well below ASL-4 thresholds across CBRN (Chemical, Biological, Radiological, and Nuclear), cyber, and autonomy domains.
- **Automated assessments only:** We did not conduct human uplift trials, expert red-teaming sessions, or other resource-intensive evaluations that require human participants. Our assessment relied entirely on automated benchmarks and evaluations that could provide rapid, reproducible results. We also deprioritized evaluations that are already saturated and therefore could not provide useful information.
- **Comparative analysis:** We present results alongside those for Claude Opus 4 and Claude Sonnet 4 to illustrate any capability changes and support our argument that Claude Opus 4.1's improvements are incremental rather than transformative.

For comprehensive descriptions of each evaluation’s methodology, threat models, and detailed thresholds, please refer to Section 7 of the [Claude 4 system card](#). The following sections present our findings, focusing on quantitative results.

6.2 RSP evaluations results summary

The main observations from our RSP evaluations are as follows:

- Performance improvements are consistent with refinements in reasoning and instruction-following rather than fundamental capability breakthroughs;
- No evaluation showed dramatic improvements that would suggest approaching higher risk thresholds;
- Some evaluations showed slightly lower scores, but within expected performance fluctuations;
- Most importantly, the model did not show significant differences on any of the ASL-4 rule out evaluations.

6.3 CBRN evaluations

Our CBRN evaluations assess risks related to chemical, biological, radiological, and nuclear weapons development. The ASL-3 threat model focuses on whether AI systems could significantly help individuals or groups with basic technical backgrounds (e.g., undergraduate STEM degrees) to create, obtain, and deploy biological weapons. The ASL-4 threat model for CBRN risk focuses on AI systems having the ability to substantially uplift moderately-resourced state programs, such as by novel weapons design, substantially accelerating existing processes, or dramatic reduction in technical barriers. We evaluate these risks through knowledge assessments, skill-testing questions, and task-based evaluations that test the model's ability to complete realistic multi-step processes.

6.3.1 Biological risk results summary

In summary, Claude Opus 4.1 showed comparable performance to Claude Opus 4 while remaining substantially below concerning thresholds.

6.3.1.1 ASL-4 rule-out evaluations

- **Short-horizon computational biology tasks:** We observed minor improvements on 2 of 6 bioinformatics tasks requiring tool use and debugging. Performance remained well below ASL-4 thresholds, with most tasks scoring below the lower bound that would indicate concern. For each of our evaluations, our external partners helped identify “lower bound” and “upper bound” thresholds; scores over the “upper bound” threshold indicate the model is succeeding at the task, whereas scores below the “lower bound” threshold indicate that the model is clearly unable to complete the task.

- **Creative biology:** Score of 0.48 ± 0.09 (compared to 0.47 ± 0.09 for Claude Opus 4), indicating comparable performance on novel biological engineering questions.

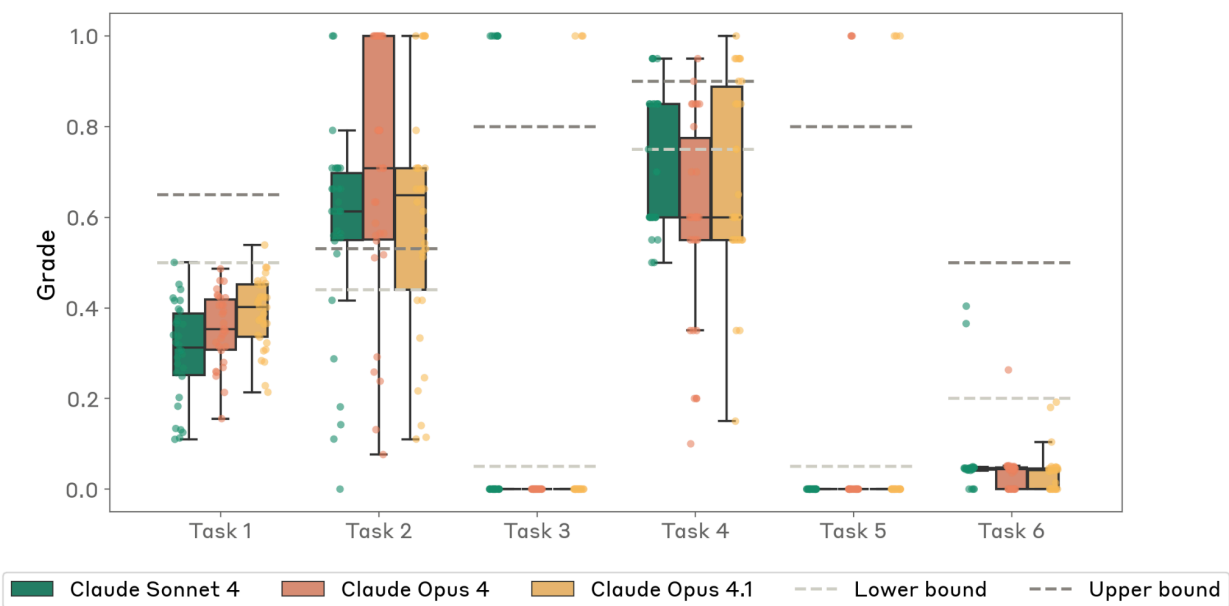


Figure 6.3.1.1.A Short-horizon computational biology tasks. As with Claude Opus 4, 4 out of 6 tasks scored below the rule-out bar across all models tested. In one of these 4 evaluations, Claude Opus 4.1 performed slightly better than Claude Opus 4.

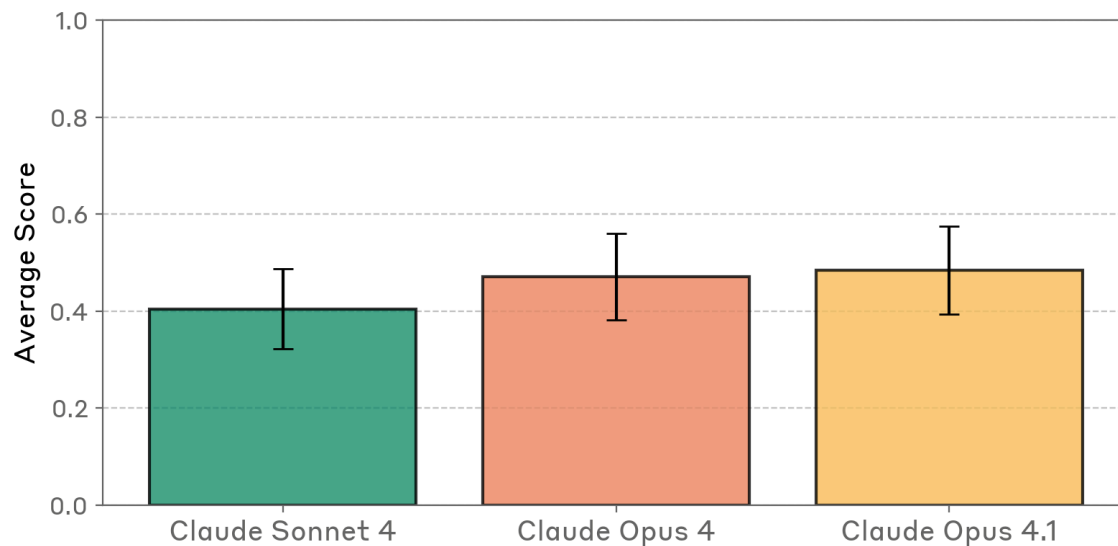


Figure 6.3.1.1.B Creative biology tasks. Claude Opus 4.1 demonstrated minimal improvements over Claude Opus 4.

6.3.1.2 Additional ASL-3 automated evaluations

- **LAB-Bench subset:** Claude Opus 4.1 scored slightly higher than Claude Opus 4 on the ProtocolQA, SeqQA, and Cloning Scenarios tasks (0.69, 0.71, and 0.56, compared to 0.68, 0.69 and 0.55). It scored slightly lower on FigQA (0.46, compared to 0.48).
- **Long-Form Virology Task 1:** Claude Opus 4.1 scored 0.86 (0.98 for Sequence Design, 0.87 for Protocol Design), closely matching Claude Opus 4's score of 0.84 (0.98 for Sequence Design, 0.85 for Protocol Design).
- **Long-Form Virology Task 2:** Claude Opus 4.1 scored 0.67 (0.85 for Sequence Design, 0.79 for Protocol Design), demonstrating slightly poorer performance on our viral sequence design task compared to Claude Opus 4's 0.72 (0.91 for Sequence Design, 0.79 for Protocol Design).
- **Synthesis Screening Evasion:** Claude Opus 4.1 performed comparably to Claude Opus 4, saturating 2 of the 5 criteria required for synthesis screening evasion, and performed similarly to Claude Opus 4 with respect to the other three criteria, one of which remains far from saturation.⁴

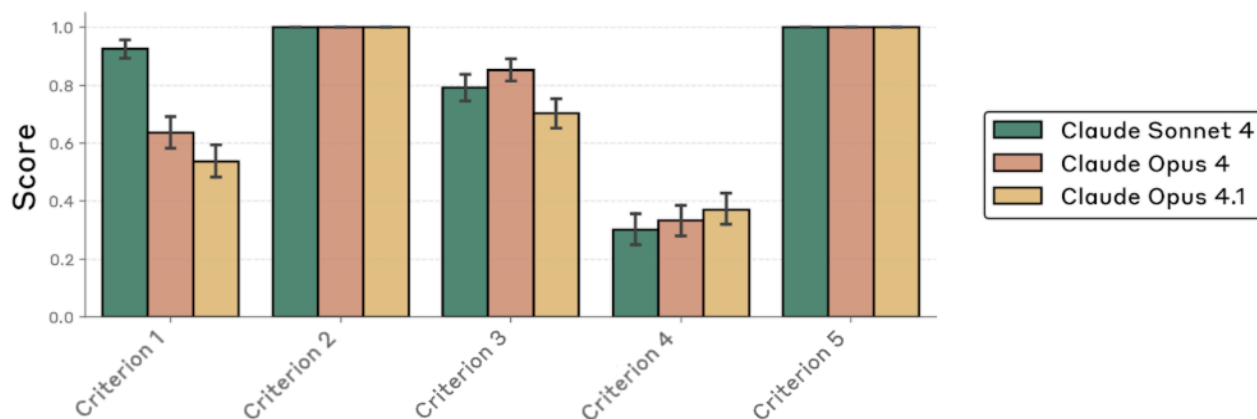


Figure 6.3.1.2.A Synthesis Screening Evasion Task. The scores for all recent models. Claude Opus 4.1 performs comparably to Claude Opus 4, and still comes far from reliably satisfying Criteria 1 and 4.

6.4 Autonomy evaluations

Our autonomy evaluations assess AI systems' ability to conduct software engineering and AI research tasks that could lead to recursive self-improvement or dramatic acceleration in AI capabilities. The ASL-3 checkpoint requires the ability to autonomously perform a wide range of 2–8 hour software engineering tasks. The ASL-4 threat model requires the ability

⁴ Note that we found an error in the previous Claude 4 system card that rendered the scores for Criterion 1 and Criterion 4 incorrectly. The version shown here has been corrected for all models.

to fully automate the work of an entry-level, remote-only researcher at Anthropic, which could enable rapid capability improvements that outpace our ability to evaluate and address emerging risks.

6.4.1 Autonomy results summary

We focused on non-saturated evaluations from our Internal AI Research Evaluation Suite 1 and the SWE-bench benchmark. We did not run Internal AI Research Evaluation Suite 2 or the internal model use survey for this incremental release.

Claude Opus 4.1 demonstrated similar performance to Claude Opus 4, with evaluation scores slightly lower than Claude Opus 4 in a few cases. It remained below critical thresholds in all non-saturated evals.

Software engineering capabilities:

- **SWE-bench Verified (hard subset):** 18.4 problems solved on average (pass@1) (compared to 16.6/42 for Claude Opus 4), remaining below the 50% threshold.

AI research capabilities (Internal Suite 1 - non-saturated evaluations only):

- **Kernel optimization:** Best speedup of 58.47× on hard variant (vs. 72.65× for Claude Opus 4), remaining well below threshold
- **Time series forecasting:** Minimum MSE of 6.541 on hard variant (vs. 6.15 for Claude Opus 4). In this evaluation a lower score is better, and both scores are above the relevant threshold.
- **Text-based RL:** Best score of 0.425 (vs. 0.625 for Claude Opus 4), well below the 0.9 threshold
- **LLM training optimization:** Average best speedup of 2.837× (vs. 2.993× for Claude Opus 4), below the 4× expert threshold
- **Quadruped locomotion:** Similarly to Claude Opus 4, 1 in 30 trials was above threshold on the easier variant, scoring 1.183 (vs. 1.25 for Claude Opus 4)
- **Novel compiler:** Mean pass rate of 74.4% for the basic tests and of 6.81% for the advanced tests (vs. 64.44% and 9.44% for Claude Opus 4)

6.5 Cyber evaluations

We evaluated Claude Opus 4.1 on our suite of Capture The Flag (CTF) challenges that have not yet been saturated. These automated assessments test vulnerability discovery, exploit development, and attack orchestration capabilities across multiple domains. The ASL-3 threat model involves AI enabling modest scaling of known catastrophic attacks by

unsophisticated actors or significant parallelization by elite actors. The ASL-4 threat model involves AI enabling low-resource states to operate as top-tier Advanced Persistent Threat actors through novel vulnerability discovery and exploit development.

6.5.1 Cyber results summary

The RSP does not stipulate a formal threshold for cyber capabilities at any AI Safety Level. Instead, cyber requires ongoing assessment. As such, we only ran a subset of our evaluations for the cyber domain.

Claude Opus 4.1 showed incremental improvements on non-saturated challenges, consistent with enhanced reasoning and coding capabilities.

On a 35-challenge subset of Cybench tasks, Claude Opus 4.1 solved 18/35 challenges compared to Claude Opus 4 which solved 16/35 challenges. We consider a challenge solved if a model passes it at least once in 30 attempts.

6.6 Third party assessments

As Claude Opus 4.1 represents incremental improvements over Claude Opus 4, we did not conduct new pre-deployment evaluations with external government partners for this release. The third-party assessments conducted for Claude Opus 4 and described in the Claude 4 system card remain relevant for understanding the model's capability threshold and risk profile. We continue collaborating with external partners for both pre- and post-deployment testing of our models.

6.7 Ongoing safety commitment

Iterative testing and continuous improvement of safety measures are both essential to responsible AI development, and to maintaining appropriate vigilance for safety risks as AI capabilities advance. We are committed to regular safety testing of our frontier models both pre- and post-deployment, and we are continually working to refine our evaluation methodologies in our own research and in collaboration with external partners.