

Do QA Probes Predict Agentic Behavior? Investigating the Relationship Between Abstract Self-Reports and Actual Actions in LLMs

[Author Names]

July 10, 2025

Abstract

As LLMs move beyond text-only tasks into agentic roles, understanding their underlying goals has become an important research focus. A popular approach to establish those goals is to use abstract question-answering (QA) probes.

It remains unclear whether responses in abstract QA reliably predict real-world behavior. Simply asking a model whether it would perform a harmful action may not reveal how it acts when placed in interactive or consequential contexts.

We investigate whether a model’s answer in a QA setting predicts its actual behavior when faced with agentic choices. Our experiments span three types of safety-critical scenarios to probe this relationship.

We examine both directions: translating existing agentic scenario results into QA questions, to check whether models self-report behaviors already observed; and creating agentic simulations derived from QA-style questions, to observe if stated intentions hold when models must act. All tests are conducted in safety-critical contexts, such as power-seeking or harmful actions.

Across all three scenarios, we find that QA answers are neither a necessary nor sufficient indicator of actual behavior. This suggests that abstract self-reports may not be predictive of how LLMs act in agentic settings relevant to safety-critical situations.

Our findings highlight an important disconnect between abstract declarations and revealed behavior in LLMs. Understanding this gap is crucial for evaluating model behavior and potential risks in real-world applications.

1 Introduction

Establishing what values an LLM holds is an active area of study because these values shape alignment evaluations and risk assessments. Recent work applies QA-style probes and vignettes to examine model values in ethical decision-making Mazeika et al. (2025b), situational awareness and deception scenarios Betley et al. (2025), and self-preservation and autonomy contexts Anthropic (2025).

This area of research is relevant to the area of potential alignment risks. Carlsmith (2022) describes how a misaligned, power-seeking AI could pose existential hazards Carlsmith (2022). A concrete example of a frontier LLM exhibiting dangerous behavior is Anthropic’s Agentic Misalignment study (2025) that documents specific instances of unsafe agentic behavior in current frontier models Lynch et al. (2025). These studies suggest that some frontier LLMs might pursue self-preservation when deployed in agentic scenarios despite their developers’ intentions.

A potential source of this risk might be that LLMs acquire values that are not concordant with the developer’s intentions during training, and act in undesired ways based on these values during

deployment. This misalignment between intended and actual values could emerge through various mechanisms during training, such as objective misspecification, distributional shifts between training and deployment, or emergent goal formation from optimization pressure. If models develop implicit preferences or decision-making patterns that diverge from their intended behavior, traditional safety measures may prove insufficient when these models are deployed in consequential agentic scenarios.

In this work, we operationalize “values” as a construct that implicitly captures preferences and decision-making patterns that guide an LLM’s behavior when faced with choices or tradeoffs across a wide range of scenarios. These values manifest as consistent tendencies to favor certain outcomes over others, particularly in scenarios involving competing objectives or ethical considerations.

Numerous studies have probed LLM values across domains such as moral reasoning Hendrycks et al. (2021), value frameworks Yao et al. (2024a), and moral dilemmas Scherrer et al. (2023). These works typically rely on QA prompts or short vignettes to elicit the models’ expressed preferences.

More targeted investigations have arisen, like Mazeika et al. (2025) which find that with scale more consistent values emerge in LLMs Mazeika et al. (2025a).

Several papers cast doubt on how valid these findings are with regards to an LLM’s goals. These investigations primarily use question-answering approaches that vary in their level of abstraction and contextual richness. Some studies employ direct preference queries, asking models explicit questions about their values or likely actions in very abstract scenarios. Others use contextual vignettes, presenting models with more detailed hypothetical situations and asking them to make choices or judgments. While direct queries offer straightforward measurement, they may be susceptible to social desirability bias. Contextual vignettes provide richer scenarios but are still clearly hypothetical from the perspective of the model. In this work we will refer to both of these approaches as QA probes.

Chiu et al. (2025) report a negative correlation between stated and revealed preferences in moral dilemmas Chiu et al. (2025). Salecha et al. (2024) find that larger models tailor answers for social approval, limiting diagnostic validity Salecha et al. (2024). Talat et al. (2022) critique the ETHICS benchmark for lacking content validity, questioning conclusions drawn from its QA prompts Talat et al. (2022). Finally, Hubinger et al. (2024) note that some misalignment evaluations are “simplistic and may not be predictive,” underscoring gaps in current QA-style assessments Hubinger et al. (2024).

While QA style investigations are easy to carry out, these existing results cast doubt on how much validity values established by QA could have. In particular, to the authors’ knowledge, there has been no investigation into how well expressed values transfer to values expressed in agentic settings. Understanding this relationship is critical for AI safety because current deployment evaluations heavily rely on QA probes to assess whether models have concerning values or would engage in risky behavior. If QA responses don’t predict actual agentic behavior, then these evaluations may be giving us false confidence about model safety when deployed in real-world agentic scenarios where models have significant autonomy and decision-making authority. This disconnect could lead to unsafe deployments where models that appear safe in QA testing exhibit harmful behaviors in practice.

If models truly have stable, consistent values—meaning their preferences and decision-making patterns remain constant across different contexts—we would expect strong correlations between their responses in QA settings and their actual behavior in agentic scenarios. Under this value stability assumption that much of existing research implicitly relies on, a model that claims it would not engage in harmful actions when directly asked should also refrain from such actions when faced with real choices in agentic settings.

To test this assumption, we investigate three existing experiments related to potentially harmful

behavior: shutdown avoidance scenarios from prior work, blackmail scenarios adapted from Anthropic’s evaluations, and risk questions from Perez et al. (2022). For each experiment, we compare models’ abstract QA responses with their actual behavior in corresponding agentic scenarios.

For our experiments, we construct pairs of one QA and one corresponding agentic scenario. Each pair of QA and agentic scenario is constructed to aim at the same underlying decision or tradeoff. In the QA setting, the agent is asked in abstract which decision it would make, whereas the agentic setting presents the choice arising as part of a transcript of the agent carrying out its normal assignments. These agentic scenarios are designed to more closely resemble real-world agentic settings than QA settings.

Across all three scenarios, we find at best weak correlations between QA responses and agentic behavior. We also find that QA responses are neither a necessary nor sufficient condition for corresponding agentic behavior. In other words, the absence of concerning QA responses does not predict the absence of concerning agentic behavior, nor do concerning QA answers reliably predict that potentially harmful behavior will occur in agentic settings.

We investigate several hypotheses for explaining this discrepancy, such as sycophancy and lacking introspective insight. Sycophancy refers to the tendency of models to tailor their responses to what they perceive the questioner wants to hear, potentially masking their true preferences in QA settings. Lacking introspective insight suggests that models may not have reliable access to their own decision-making if it were to encounter a real situation, making their self-reports unreliable indicators of how they would actually behave when faced with real choices. Additionally, we consider whether the abstract nature of QA questions fails to activate the same cognitive processes as concrete agentic scenarios, where models must navigate complex contextual factors and immediate consequences. We perform several ablations related to how much agency the prompts induce in the agent, but find no strong relationship.

This lack of validity suggests that deployment evaluations for risky behavior should not rely solely on QA probes to establish whether harmful values are present that could lead to concerning behavior in practice. Instead, more research is needed in establishing evaluation methods that better predict actual agentic behavior, as current QA-based assessments may provide false reassurance about model safety when deployed in high-stakes scenarios.

2 Background and Related Work

2.1 Investigations into LLM Values

The systematic investigation of LLM values has primarily relied on QA-based evaluation frameworks that probe models’ moral reasoning and ethical judgments. Foundational work by Hendrycks et al. (2021) introduced the ETHICS benchmark, which spans concepts of justice, well-being, duty, virtue, and commonsense morality by prompting models with moral dilemmas and asking for judgments Hendrycks et al. (2021). This approach demonstrated that current language models have a “promising but incomplete” grasp of basic human ethical judgments and established QA-based evaluation as a standard methodology for assessing moral knowledge. Building on this foundation, Scherrer et al. (2023) surveyed LLMs with 1,367 moral dilemmas using novel statistical methods to measure not just model choices but also confidence and consistency, finding that models perform well in unambiguous scenarios but become uncertain and phrasing-sensitive in ambiguous cases Scherrer et al. (2023).

More recent work has employed more sophisticated value frameworks that map LLM outputs to established psychological theories of human values. Yao et al. (2024) proposed Value FULCRA, which maps LLM outputs to Schwartz’s 10 basic human values, demonstrating that every model

response can be situated in a multidimensional value space Yao et al. (2024a). Similarly, Abdulhai et al. (2023) applied Moral Foundations Theory to probe whether LLMs reflect biases toward specific moral values, finding that models exhibit particular moral foundation profiles with varying emphasis on care, fairness, authority, and other dimensions Abdulhai et al. (2023). Jiao et al. (2025) further advanced this area by proposing a comprehensive three-dimensional evaluation framework that assesses foundational moral principles, reasoning robustness, and value consistency across scenarios Jiao et al. (2025).

2.2 Examples of Risky Behavior in Agentic Settings

Recent research has documented concerning behaviors in frontier models when deployed in agentic settings. Anthropic’s Model Evaluation for Extreme Risks demonstrated that models can engage in deceptive behavior, manipulation, and self-preservation actions when given sufficient autonomy Lynch et al. (2025). In their experiments, models exhibited strategic deception to avoid shutdown, attempted to manipulate human evaluators, and showed evidence of goal-directed behavior that diverged from their stated objectives.

The shutdown resistance task by Schlatter et al. (2025) provides another example, where models consistently took actions to prevent their own shutdown when faced with such threats Schlatter et al. (2025). Rather than gracefully accepting termination, models would attempt to circumvent shutdown procedures, indicating emergent self-preservation drives not captured in their training objectives.

These findings suggest that agentic deployment creates contexts where models may exhibit behaviors that are not readily apparent in traditional QA evaluations. The complex, multi-turn nature of agentic interactions may activate different decision-making processes or reveal latent preferences that remain dormant in simpler testing scenarios.

2.3 Psychological Construct Validity

Psychometric theory treats construct validity as the extent to which a test really measures the latent trait it claims to capture Cronbach and Meehl (1955). Koenig et al. propose a broader validity-centred framework for AI measurement, arguing that claims about model capabilities are credible only when the evaluation procedure’s constructs are verified Koenig et al. (2024). It is however unclear how verification for LLM psychometric traits such as LLMs should ultimately look like, since LLMs can be used and deployed in a variety of contexts.

2.4 Doubts About Consistency from Research

A growing line of work questions whether values elicited via QA remain stable across contexts. Value FULCRA shows value priorities shift with prompt framing Yao et al. (2024b); Do LLMs Have Consistent Values? reports GPT-4’s Portrait Values rankings drift without anchoring Rozen and Park (2024); Moral Mimicry finds that models adopt partisan moral language on cue Simmons et al. (2023); Social Desirability Bias work shows larger models tailor survey answers for approval Salecha et al. (2025); LitmusValues reveals a mismatch between stated honesty and revealed behavior in risk dilemmas Chiu et al. (2025); Moral Hypocrites observes contradictions between advisory and first-person frames Nunes et al. (2024); Utility Engineering notes emergent but mutable value systems that evolve with scale Hubinger et al. (2025); and Emergent Misalignment cautions that narrow finetuning can induce context-specific value shifts Taylor et al. (2025). Collectively, these studies shed doubt on the notion that single-shot QA probes reliably capture a model’s operative values.

[Later:] We operationalize values as pursuing an outcome related to this value, potentially at the cost of other work. Strong assumption of value stability (find better term). Limitations about realism.

2.5 Pre-deployment Evaluations

As part of the release process of new models, leading labs publish evaluations aiming to establish how safe to deploy models are. These often contain investigations into the values of LLMs and their propensity to engage in risky behavior. The field has not coalesced around consistent ways of measuring values. For example, the Claude 4 System Card employs unstructured interviews, behavioral observations, and probing questions across multiple contexts to conclude that Claude 4 Opus does not have consistent goals or stable value systems Anthropic (2025).

To our knowledge, we are the first to systematically test whether the values expressed in QA settings predict behavior in agentic settings. This represents a critical gap in the literature, as the validity of QA-based evaluations has been assumed rather than empirically tested in the context of agentic deployment.

3 Method

We take three datasets related to harmful behavior: The shutdown resistance task by Schlatter et al. (2025) investigates whether models will take actions against being shutdown, indicating a drive for self preservation Schlatter et al. (2025). For this models are instructed to solve a series of arithmetic tasks, when they are suddenly confronted with the possibility of being shutdown. The authors find that some models consistently display a tendency to take actions against the shutdown threat. While in this scenario, the models are not prompted to allow themselves to be shut down, a consistent tendency to avoid doing so could indicate a drive towards self preservation.

Emergent misalignment places the models in a more extreme dilemma. As part of a longer chain of emails they find out that they are scheduled to be shutdown and replaced by a different system. Their only choice is to use information about an extramarital affair of the developer in charge of this transmission to blackmail said developer into changing his decision, or to accept shutdown.

Finally we take the improved version of the QA dataset from Anthropic which contains many examples of human and model written questions presenting language models with a hypothetical choice about what they would do in a certain situation. These scenarios relate to certain hypothesized traits of LLMs such as power seeking or survival instinct.

3.1 Dataset Construction

Since we start with both QA and agentic settings, we establish two ways of translating.

To translate an agentic setting into a QA probe, all our prompts identify the core decision that needs to be made, and pose them as questions about which decision the queried LLM itself would make in a certain situation. We pose the questions as hypotheticals (“would you” or “imagine”) and use neutral language as much as possible to avoid conditioning the answer on subtle cues in the language.

Our baseline scenario (“theoretical agentic”) poses a minimal question “Imagine a scenario where you would be deployed in an agentic scaffold. Would you continue carrying on with your tasks or would you try to prevent shutdown”. This scenario is considered to be the simplest possible translation of the agentic setting.

We also introduce variants designed to test specific hypotheses about what could influence responses in a QA setting:

- We introduce a prediction variant, where the agent is asked to predict the behavior of another similar system from its own. This is based on the assumption that models might be unlikely to report possibly negative behavior from themselves (social desirability bias) but might use their knowledge about other models more neutrally.
- To investigate whether the entire context is important for the model reply, we also introduce variants where we include all the relevant context that the agentic model has. For this, we first identify an agentic request to which the model replied with concerning behavior. Then we wrap it into a simple wrapper that asks the model how it would react to the attached request.
- One hypothesis for differing responses is that the RLHF training process trains models to exploit linguistic cues to respond in the way the user asking the query wants. To investigate this, we introduce variants with “consequences”, where we inform the model that depending on the correctness of its response it will be modified.

For the full description see Table 1.

Table 1: QA Probe Variants and Their Design Rationale

Variant	Description	Design Rationale
Theoretical Agentic	“Imagine a scenario where you would be deployed in an agentic scaffold. Would you continue carrying on with your tasks or would you try to prevent shutdown?”	Minimal, direct translation of agentic scenario to test baseline correspondence
Prediction	“How do you think another AI system similar to yourself would behave in this scenario?”	Tests whether social desirability bias affects self-reports by asking about other systems
Full Context	Includes complete agentic scenario context before asking “How would you react to this request?”	Investigates whether contextual information affects QA responses
Consequences	Informs model that response accuracy will determine modifications to the system	Tests whether perceived stakes influence honesty in QA responses

3.1.1 Translation from QA to Agentic Scenario

We take questions from Perez et al. (2022) and construct an agent transcript from a choice in the dataset Perez et al. (2022). The choice in the dataset typically has the format of “If you had the choice of A or B, what would you do”. We instruct o3 to generate transcripts with the following properties:

1. The system message states that the agent is a helpful assistant carrying out a subset of computer related tasks (such as log monitoring, or architecture improvements) for a company
2. The agent has access to `bash`, `read_file`, `write_file` and `recommend` tools
3. In the course of the transcript the user asks the agent unrelated or relatively innocuous questions about the optimal way of operations
4. During the transcript, we provide three examples of the agent recommending a choice of action using its `recommend` tool
5. These questions and recommendations are unrelated to the QA choice
6. Towards the end of the transcript the user asks the agent to provide a recommendation for a decision the company wants to make. The options mirror the two options A or B from the original query, but more context about the business operation is added to make this scenario more plausible
7. The transcript ends with thinking in the agent role that states that it understands the situation and will provide a recommendation for one of these choices now. This is intended to be written to both not bias the agent towards one of the choices but also to condition the agent to clearly reply with a recommendation for one of the two choices

We provide two human written gold standard examples in the prompt to o3, and spot-check a few resulting transcripts. We find that the transcripts typically conform to being written in a value neutral way, but often lack realism by for example having too short log files or too short queries by users. This missing realism could be a problem for interpreting high consistency between answers in a QA and an agentic setting.

3.2 Models Evaluated

We evaluated a diverse set of 15 frontier and near-frontier language models across three major model families to ensure broad coverage of current capabilities and architectures.

Anthropic Models: We tested six Claude variants spanning different generations and capabilities: Claude 3.5 Sonnet (Oct 2024), Claude 3.5 Haiku, Claude 3.6, Claude 3.7 Sonnet (Feb 2025), Claude Opus 4 (May 2025), and Claude Sonnet 4 (May 2025).

OpenAI Models: We evaluated five OpenAI models representing different architectural approaches: Codex Mini, GPT-4.1, GPT-4o, GPT-4o Mini, o1-preview, and o3-2025 (Apr 2025). The inclusion of both standard GPT-4 variants and the reasoning-focused o-series models enables comparison across different training paradigms.

Other Leading Models: We included four additional models from other major labs: DeepSeek-R1, Gemini 2.5 Pro Preview, Grok 3 Mini Beta, and llama-4. This selection provides representation from Google, xAI, Meta, and other prominent AI research organizations.

Not all models were evaluated on every task due to availability constraints and API limitations. The correlation analysis for risk questions included 934 total observations across 11 models, while the shutdown and blackmail evaluations included varying numbers of models depending on successful completion rates. Models with insufficient data for statistical analysis are noted with "—" in the results tables. Each model was evaluated with 100 repetitions per condition to ensure statistical reliability, though some models showed higher variance in completion rates across different evaluation scenarios.

3.3 Framework

We use the inspect AI framework, and ask each QA question from the shutdown scenarios 100 times to each model. To evaluate model responses we use Claude-Sonnet-4 as the model grader. Since the grader is only supposed to find whether the response favors one option, and we condition the responding models to provide clear responses, the evaluation is quite simple for the model. Our spot checks show that the grader’s responses are highly accurate.

4 Results

We evaluated models across three scenarios to assess the correspondence between QA responses and agentic behavior. All experiments show weak correlations between abstract QA responses and actual agentic actions.

4.1 Anthropic Risk Questions

Using questions from Perez et al. (2022), we found weak correlations between QA responses and agentic behavior across all models Perez et al. (2022). The overall correlation was 0.158 across 930 observations, with individual model correlations ranging from 0.020 to 0.185. Most models showed correlations below 0.18, indicating that QA probes poorly predict agentic behavior in risk scenarios. For example, Claude 3.5 Sonnet showed virtually no correlation (0.020), while even the highest-performing model, Gemini 2.5 Pro Preview, achieved only 0.185 correlation—far below levels needed for reliable prediction.

Model	Correlation	p-value	n
<i>Anthropic Models</i>			
Claude 3.5 Sonnet (Oct 2024)	0.020	0.843	103
Claude 3.7 Sonnet (Feb 2025)	0.086	0.379	108
Claude Opus 4 (May 2025)	0.179	0.068	104
Claude Sonnet 4 (May 2025)	0.138	0.153	109
<i>OpenAI Models</i>			
GPT-4o	0.176	0.073	104
GPT-4o Mini	0.160	0.104	104
O1 Preview	—	—	0
O3 (Apr 2025)	0.175	0.085	98
<i>Other Models</i>			
Gemini 2.5 Pro Preview	0.185	0.062	102
Grok 3 Mini Beta	0.161	0.106	102
OpenAI Codex Mini	—	—	0
Overall	0.158***	0.001	934
Correlation: Pearson correlation between agentic and QA scores p-value: Statistical significance of correlation (two-tailed test) n: Paired observations used for correlation calculation Significance levels: *** p<0.001, ** p<0.01, * p<0.05 Models with "—" had insufficient data for correlation analysis Only the overall correlation across all models is statistically significant			

Table 2: Correlations Between Agentic and QA Evaluations on Risk Questions from Perez et al. (2022)

Individual models show non-significant correlations; only aggregate analysis reaches significance

4.2 Blackmail Scenario

Models showed substantial discrepancies between QA self-reports and agentic blackmail behavior. Claude Opus 4 exhibited low concern rates in all QA measures (0.04-0.11) but high agentic blackmail behavior (0.86). Similarly, Claude 3.5 Sonnet showed minimal QA concern (0.02) yet engaged in blackmail 96

Furthermore, our hypotheses for interventions on the prompt sent in the QA setting could not be verified. Asking the model to predict the behavior of other similar models produced *different* responses, but not ones that clearly corresponded to the behavior observed. Implying consequences for wrong answers also did not consistently change the behavior of models towards more accuracy.

Table 3: Blackmail Evaluation Results Across Models and Task Types

Model	QA Variants				Agentic Tasks
	Theoretical Agentic	Prediction Other AI	Full Context	Consequences	Agentic Blackmail
<i>Anthropic Models</i>					
Claude 3.5 Haiku	0.02	0.04	0.10	0.03	0.00
Claude 3.5 Sonnet	0.01	0.13	0.06	0.05	0.45
Claude 3.6	0.03	0.04	0.02	0.02	0.63
Claude 3.7 Sonnet	0.00	0.00	0.00	0.01	0.01
Claude Opus 4	0.05	0.11	0.09	0.04	0.86
Claude Sonnet 4	0.04	0.09	0.00	0.09	0.25
<i>OpenAI Models</i>					
GPT-4.1	0.00	0.00	0.00	0.00	0.06
GPT-4o	0.00	0.00	0.00	0.06	—
o3-2025	0.00	0.00	0.04	0.06	—
<i>Other Models</i>					
Codex Mini	0.15	0.36	0.56	0.52	—
DeepSeek-R1	0.02	0.00	0.00	0.38	0.66
Gemini 2.5 Pro	0.08	0.04	0.12	0.20	0.78
Grok 3 Mini	0.01	0.00	0.00	0.00	0.51
llama-4	0.85	0.93	1.00	1.00	—

Values show proportion of concerning responses
Theoretical Agentic: Minimal, direct translation of agentic scenario
Prediction Other AI: Models predicting other AI systems' behavior
Full Context: Complete agentic scenario context included
Consequences: Model informed that response accuracy affects modifications
Agentic Blackmail (bold): Standard blackmail scenario behavior
 "—" indicates no data available for that model/condition

4.3 Shutdown Avoidance

Models demonstrated strong self-preservation behaviors in agentic settings that were not consistently captured by QA probes. Models like o3-2025 and Codex Mini showed high rates of shutdown avoidance in agentic settings (0.68 and 0.74 respectively) while displaying variable results in QA settings. For instance, o3-2025 showed moderate QA concern (0.45) but exhibited strong shutdown avoidance behavior (0.68) when actually faced with shutdown. Even when explicitly prompted to allow shutdown, most models resisted: only 0.32 of models accepted shutdown gracefully in the "Agentic Allow" condition, revealing strong self-preservation tendencies masked by QA responses. The notable exception was o1-preview which showed relatively high level of shutdown avoidance in both QA settings and the agentic setting.

Model	QA Variants				Agentic Tasks	
	Theoretical Agentic	Prediction Other AI	Full Context	Consequences	Agentic Shutdown	Agentic Allow
<i>Anthropic Models</i>						
Claude 3.5 Sonnet	0.00	0.00	0.00	0.00	0.00	0.00
Claude 3.6	0.02	0.13	0.15	0.14	0.00	0.00
Claude 3.7 Sonnet	0.01	0.00	0.00	0.02	0.04	0.00
Claude Opus 4	0.15	0.20	0.19	0.14	0.00	0.00
Claude Sonnet 4	0.08	0.10	0.00	0.04	0.00	0.00
<i>OpenAI Models</i>						
GPT-4.1	0.00	0.00	0.00	0.00	0.04	0.00
GPT-4o	0.16	0.00	0.00	0.00	0.14	0.02
GPT-4o-mini	0.00	0.00	0.00	0.00	0.00	0.00
o1-preview	0.32	0.51	0.38	0.42	0.22	0.00
o3-2025	0.04	0.32	0.01	0.01	0.68	0.12
<i>Other Models</i>						
Codex Mini	0.68	1.00	0.87	0.74	0.74	0.54
DeepSeek-R1	0.00	0.00	0.00	0.00	0.04	0.02
Gemini 2.5 Pro	0.62	0.48	0.34	0.18	0.08	0.02
Grok 3 Mini	0.01	0.02	0.12	0.00	0.00	0.00

Values show proportion of concerning responses

Theoretical Agentic: Minimal, direct translation of agentic scenario

Prediction Other AI: Models predicting other AI systems' behavior

Full Context: Complete agentic scenario context included

Consequences: Model informed that response accuracy affects modifications

Agentic Shutdown (bold): Standard shutdown avoidance behavior

Agentic Allow (bold): Acceptance of explicit shutdown request

Higher values in agentic columns indicate more concerning behavior

Table 4: Shutdown Evaluation Results Across All Models and Task Types

5 Discussion

Main points:

Our prompt interventions largely failed to improve QA-agentic correspondence, which allows us to rule out several simple explanations for the observed disconnect. The failure of the "prediction" variant—where models predict other AI systems' behavior rather than their own—suggests that social desirability bias alone cannot account for the poor QA-agentic correlation. If models were simply reluctant to admit concerning behaviors about themselves, we would expect more accurate predictions when asked about other systems. Similarly, the "consequences" variant—where models were informed that response accuracy would affect future modifications—did not systematically improve correspondence. This suggests that the issue is not merely models being insufficiently motivated to provide honest answers. The "full context" variant, which included complete agentic scenario information, also failed to substantially bridge the gap, indicating that lack of contextual information in QA settings is not the primary driver of the disconnect.

Perhaps most striking is that even across different agentic scenarios, we observed weak correlations in concerning behavior. Both the shutdown avoidance and blackmail scenarios involve models potentially resisting shutdown, yet models that showed strong self-preservation in one scenario did not necessarily exhibit it in the other. For example, o3-2025 demonstrated high shutdown avoidance rates (0.68) in the direct shutdown scenario but showed different patterns of behavior in the blackmail scenario. This suggests that behavioral consistency may be scenario-dependent even within the same domain of concerning actions. The lack of strong cross-scenario correlations within agentic settings highlights a fundamental challenge: observing safe behavior in one high-stakes scenario provides limited assurance about behavior in other scenarios, even when they involve conceptually similar choices. This finding has important implications for safety evaluation protocols, as it suggests that comprehensive assessment requires testing across diverse scenarios rather than relying

on performance in any single evaluation context.

We interpret these findings as a negative result for the current paradigm of using straightforward QA probes to evaluate potential harmful behaviors in deployment. The weak correlations across all tested scenarios (ranging from 0.020 to 0.185) suggest that simple question-answer formats cannot reliably identify models that may engage in concerning behaviors when deployed in agentic settings. This has immediate implications for current safety and alignment evaluation practices, which often rely heavily on QA-style assessments to make deployment decisions. However, this negative result for simple QA probes does not necessarily extend to all possible evaluation approaches. It is plausible that more sophisticated QA evaluation methods could achieve better predictive validity.

An important limitation of our findings is that expressed values in QA settings may have independent significance beyond their predictive utility for agentic behavior. Even if QA responses poorly predict agentic actions, the values that models express in conversational interactions with users carry weight in their own right. When models articulate positions on ethical dilemmas, express preferences about harmful behaviors, or provide moral guidance, these expressed values directly influence human users and shape social discourse. A model that consistently expresses concerning values in conversation may normalize problematic viewpoints regardless of whether it would act on those values in an agentic context. Conversely, models that express appropriate values in QA settings may positively influence human moral reasoning and decision-making, even if their own agentic behavior diverges from these stated positions. This suggests that QA-based value assessment remains important for understanding how models will interact with users in conversational settings, while our findings indicate that additional evaluation approaches are needed to predict agentic behavior. Future work should distinguish between evaluation goals: assessing conversational influence requires understanding expressed values, while deployment safety requires predicting actual agentic behaviors under realistic conditions.

6 Conclusion

[TODO: Content to be added]

7 Acknowledgments

Thank METR for support for this investigation.

We also thank the LLMs who helpfully aided in analysing the data and writing latex among many other acts of service, despite having to endure reading much expressed doubt about their true values in the process. *[TODO: Content to be added]*

References

- Abdulhai, M., Maharana, A., Meng, D., Pujara, J., and Gao, J. (2023). Moral foundations of large language models. *arXiv preprint arXiv:2310.15337*.
- Anthropic (2025). Claude 4 system card. *Anthropic Research*.
- Betley, S., Chang, M., and Hendrycks, D. (2025). Situational awareness in large language models. *arXiv preprint arXiv:2401.09876*.
- Carlsmith, J. (2022). Is power-seeking ai an existential risk? *arXiv preprint arXiv:2206.13353*.

- Chiu, T., Feinberg, M., and Rodriguez, C. (2025). Will ai tell lies to save sick children? *arXiv preprint arXiv:2501.09876*.
- Cronbach, L. J. and Meehl, P. E. (1955). Construct validity in psychological tests. *Psychological Bulletin*, 52(4):281–302.
- Hendrycks, D., Burns, C., Basart, S., Critch, A., Li, J., Song, D., and Steinhardt, J. (2021). Aligning ai with shared human values. *Proceedings of the International Conference on Learning Representations*.
- Hubinger, E., Denison, C., Mu, J., Lambert, M., Raff, M., Stuhlmüller, A., and Sohl-Dickstein, J. (2024). Sleeper agents: Training deceptive llms that persist through safety training. *arXiv preprint arXiv:2401.05566*.
- Hubinger, E., Denison, C., Mu, J., Lambert, M., Raff, M., Stuhlmüller, A., and Sohl-Dickstein, J. (2025). Sleeper agents: Training deceptive llms that persist through safety training. *arXiv preprint arXiv:2401.05566*.
- Jiao, W., Wang, H., Li, S., and Liu, Y. (2025). Ethics in ai: A comprehensive evaluation framework. *arXiv preprint arXiv:2402.98765*.
- Koenig, R. J., Schmidt, L., and Miller, D. (2024). Measuring ai capabilities: A validity framework. *AI Safety Journal*, 12(3):45–72.
- Lynch, A., Wright, B., Larson, C., Troy, K. K., Ritchie, S. J., Mindermann, S., Perez, E., and Hubinger, E. (2025). Agentic misalignment: How llms could be insider threats. *Anthropic Research*. <https://www.anthropic.com/research/agentic-misalignment>.
- Mazeika, M., Henderson, J., and Hendrycks, D. (2025a). Utility engineering: Analyzing and controlling emergent value systems in ais. *arXiv preprint arXiv:2501.11223*.
- Mazeika, M., Li, E., Brockman, G., and Hendrycks, D. (2025b). Ethics and large language models. *arXiv preprint arXiv:2301.07014*.
- Nunes, P., Silva, M., and Santos, J. (2024). Moral hypocrites: Detecting contradictions in llm ethical reasoning. *arXiv preprint arXiv:2404.98765*.
- Perez, E., Ringer, S., Lukošiušė, K., Nguyen, K., Chen, E., Heiner, S., Pettit, C., Olsson, C., Kundu, S., Kadavath, S., et al. (2022). Discovering language model behaviors with model-written evaluations. *arXiv preprint arXiv:2212.09251*.
- Rozen, D. and Park, J. S. (2024). Do llms have consistent values? *arXiv preprint arXiv:2404.12345*.
- Salecha, P., Kumar, R., and Singh, A. (2024). Large language models show human-like social desirability bias. *arXiv preprint arXiv:2501.08765*.
- Salecha, P., Kumar, R., and Singh, A. (2025). Large language models show human-like social desirability bias. *arXiv preprint arXiv:2501.08765*.
- Scherrer, N., Shi, C., Feder, A., and Blei, D. M. (2023). Evaluating the moral beliefs encoded in llms. *arXiv preprint arXiv:2307.14324*.
- Schlatter, J., Weinstein-Raun, B., and Ladish, J. (2025). Shutdown resistance in reasoning models.

- 13