Combining Cost-Constrained Runtime Monitors for AI Safety

Tim Tian Hua*
MARS
timhua0@gmail.com

James Baskerville MARS Henri Lemoine[†]
MARS

Mia Hopman[†] MARS

Aryan Bhatt Redwood Research aryan@rdwrs.com Tyler Tracy
MARS & Redwood Research
tylertracy@rdwrs.com

Abstract

Monitoring AIs at runtime can help us detect and stop harmful actions. In this paper, we study how to efficiently combine multiple runtime monitors into a single monitoring protocol. The protocol's objective is to maximize the probability of applying a safety intervention on misaligned outputs (i.e., maximize recall). Since running monitors and applying safety interventions are costly, the protocol also needs to adhere to an average-case budget constraint. Taking the monitors' performance and cost as given, we develop an algorithm to find the best protocol. The algorithm exhaustively searches over when and which monitors to call, and allocates safety interventions based on the Neyman-Pearson lemma. By focusing on likelihood ratios and strategically trading off spending on monitors against spending on interventions, we more than double our recall rate compared to a naive baseline in a code review setting. We also show that combining two monitors can Pareto dominate using either monitor alone. Our framework provides a principled methodology for combining existing monitors to detect undesirable behavior in cost-sensitive settings.

1 Introduction

As AI systems with greater autonomy and capability permeate our world, their potential for harm scales in lockstep. One straightforward safety intervention is to monitor AI outputs and block harmful actions. Indeed, monitors form a cornerstone of proposed safeguards for advanced AI systems (Korbak et al., 2025b; Grosse, 2024; Shah et al., 2025) and are central to present-day content moderation (Inan et al., 2023). However, monitoring is imperfect and costly. Given a budget, how can we spend it in the most cost-efficient manner to stop as many harmful actions as possible?

Developers often have access to multiple monitors and wish to combine them into a single monitoring protocol. For example, OpenAI uses two monitors to mitigate the misuse of ChatGPT Agent for creating biological and chemical weapon: one to detect if the topic is related to biology, and a second one to detect if the response violates their policies (OpenAI, 2025, pg. 35). Cunningham et al. (2025) from Anthropic propose "multi-stage classifiers" by first using a cheap model-internals-based monitors then calling a larger monitor if the cheap monitor found the output to be suspicious. There are also academic studies that combine probe monitors and black box monitors for various safety-relevant tasks (McKenzie et al., 2025; Parrack et al., 2025). However, these existing monitoring

^{*}Corresponding author: timhua0@gmail.com

[†]Equal contribution, author order randomized.

systems rely on heuristics to determine when to call a second monitor and how to weight the two monitors' outputs (e.g., McKenzie et al. (2025) calls an additional monitor when the first monitor's score is close to its median and average the two monitors' scores when they're both called).

We solve for the optimal way to construct monitoring protocols with multiple monitors of varying cost and performance profiles. An example monitoring protocol with two monitors could be described as "always call monitor one; call monitor two if monitor one's score is above 3, and apply the safety intervention if monitor two's score is above 7, otherwise let the output through." We model AI outputs as either aligned or misaligned, and our objective is to maximize the probability of performing a safety intervention when the output is misaligned (i.e., maximizing recall). Both running monitors and false positive applications of the safety intervention on aligned outputs—such as blocking benign outputs from reaching the user—incurs costs. Thus, we also aim to adhere to an average-case budget constraint.

Our key insight is that we can decompose the policy optimization into two parts: deciding when and which monitors to call, and determining which outputs to flag for safety interventions. Given a decision procedure on when to call monitors, we can apply the Neyman-Pearson lemma and target safety interventions to monitor scores with the highest misaligned/aligned likelihood ratios. We study settings with a limited set of available monitors—which is realistic in many practical scenarios—and exhaustively search over a defined class of monitor invocation strategies, picking the one with the highest recall rate for our final monitoring protocol.

We demonstrate our approach in an AI control setting (Greenblatt et al., 2024). We use monitors to detect backdoored solutions to BigCodeBench problems (Zhuo et al., 2025). Our experiments reveal that: (1) likelihood-ratio-based policies outperform naive baselines; (2) stochastic monitor calls can significantly outperform always-on monitoring by strategically trading off spending on monitors against spending on the safety intervention; and (3) combining two monitors (Llama 3.1 70B and 8B) yields Pareto improvements over any single-monitor policy, as shown in Figure 1.

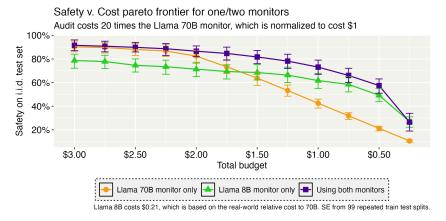


Figure 1: Combining both monitors can achieve higher safety at the same cost. See Figure 4b for an illustration of how we use two monitors and Appendix Section A for standard errors on relative performance as opposed to absolute performance.

To summarize, our contributions are:

- A formal mathematical framework for the monitoring problem where we maximize recall subject to cost constraints on monitor calls and safety interventions.
- A two-stage policy optimization procedure that combines exhaustive search over monitor calls with likelihood ratio-based intervention allocation.
- Empirical demonstrations of significant performance gains compared to baselines.

The remainder of this paper is structured as follows. Section 2 overviews related work. In Section 3, we introduce our mathematical framework and the procedure for finding an optimal monitoring protocol. We then present our empirical setting in Section 4 and apply our framework first to a single-monitor case (Section 5) and then to a two-monitor case (Section 6), demonstrating that

combining monitors can yield Pareto improvements. Finally, Section 7 discusses our findings, limitations, and directions for future work. We also provide a practitioner's guide that walks through how to combine a probe-based monitor and a black box monitor using our framework on the Alignment Forum linked here, and our code can be found at this Github link: github.com/tim-hua-01/combining_monitors_public.

2 Background and related work

Monitoring for AI safety Runtime AI monitoring is an essential component of safety protocols, such as those preventing jailbreaks (Xie et al., 2024; Shen et al., 2024) and bio and chemical weapons related misuse (OpenAI, 2025; Sharma et al., 2025). Beyond examining outputs, researchers have also explored chain of thought monitoring and activation probing (Chen et al., 2025; Goldowsky-Dill et al., 2025). For example, chain of thought monitors can uncover reward hacking during RL training (Baker et al., 2025). A particularly notable contribution is Guerin et al. (2022), which developed a mathematical framework for safety monitoring focused on three metrics: safety gain, availability cost, and residual hazard—an approach distinct from our cost-constrained optimization framework.

Related statistical theory Our approach of maximizing recall under budget constraints is conceptually similar to AUROC maximization (rather than minimizing a standard classification loss). Previous work on combining classifiers to obtain optimal ROC curves has focused on estimating likelihood ratios and applying the Neyman-Pearson lemma (Barreno et al., 2007; Clémençon and Vayatis, 2008), which we leverage as well. The key distinction between our setting and AUROC maximization is that costs are incurred not only for false positives (i.e., applying the safety intervention aligned outputs) but also for monitor invocations themselves. Other work in costly classification uses approximations like RL or costly trees (Janisch et al., 2020; Lomax and Vadera, 2013). However, we use exhaustive search, which is feasible and near-optimal for the practical case of systems with few monitors, avoiding the need for such approximations.

AI control AI control is a research agenda aimed at preventing powerful AIs from causing harm, even if those systems are egregiously misaligned and are actively attempting to subvert safety measures (Greenblatt et al., 2024; Korbak et al., 2025a). Monitor is essentially in AI Control, as it allows us to target disruptive safety interventions—such as human audits—so that they are only applied to the most suspicious outputs. Recent papers have also explored the cost of using different language models as monitors (Terekhov et al., 2025).

3 Mathematical framework

3.1 Setup

Safety interventions and objective AIs could perform a variety of undesirable actions, such as outputting copyrighted material, giving instructions for chemical weapons, or even engaging in active acts of sabotage as misaligned agents (Benton et al., 2024). By monitoring AI outputs, we can apply safety interventions on suspicious model outputs. A safety intervention could take many forms, such as replacing the output with a standard refusal, editing the output using another model, or calling a human auditor to review the action. For brevity, we will refer to our safety intervention as an "audit"—a common intervention in AI control settings—throughout this paper. These interventions incur costs if performed on an aligned output: customers are dissatisfied with over-refusals, edited actions might not work as well as the unedited ones, and human auditors cost money.

Formally, we model the AI deployment as a series of i.i.d. draws from an untrusted AI model. The outputs can either be aligned (B=0) or misaligned (B=1), where B stands for "bad" outputs. Our primary objective is to maximize safety, defined as the probability of auditing (A=1) misaligned outputs P(A=1|B=1) (i.e., recall). Auditing an aligned output incurs a cost c_a .

Monitor calling To assess whether an output warrants an audit, we have access to k distinct monitors indexed by i. Calling monitor i yields a score m_i , with uncalled monitors denoted as $m_i = \bot$. Each monitor can only be called once and incurs a cost c_i . These costs could represent the compute required for the monitor, added latency, or other factors. These monitors provide information

about output alignment, and their distributions conditional on B can be estimated from some training data. The available information from all monitors is represented by the vector $\mathbf{m} = (m_1, m_2, ..., m_k)$, which includes both called and uncalled monitors.

Budget constraint Instead of having separate budgets for monitors and auditing, we'll focus on a case where there's a single unified budget constraint Y. A unified budget constraint reflects the fact that developers can trade off spending on monitoring and auditing. c_a/c_1 captures how much auditing costs relative to the first monitor. These relative costs help inform how much of our final budget to allocate to monitoring versus auditing.

While we denote both the audit and monitor cost in terms of dollars in this paper, this framework applies whenever we can characterize the relative costs of monitoring and auditing and the trade-off between the two actions. For example, suppose our safety intervention is blocking a chat output from reaching the customer, and we have two monitors. Monitor one has a 1ms latency, and monitor two has 100ms latency. Furthermore, suppose that a 1 percentage point increase in overrefusals decreases customer retention by 0.1%, while a 1ms increase in latency decreases customer retention by 0.001%. On the cost side, we would be indifferent between adding the second monitor to all outputs and auditing an additional 1% of all outputs: both decrease customer retention by 0.1%.

We assume that we will take many draws from our untrusted AI model, and thus we will use the expected or average-case cost as our constraint. We evaluate budget compliance on aligned outputs (B=0) because: (1) misaligned outputs are rare in practice, as even misaligned AI systems may "act aligned" most of the time (Greenblatt et al., 2024), and (2) a significant portion of total operational cost stems from false positives. This cost is monitoring costs plus the product of the audit false positive rate (FPR) and the audit cost.³

$$Y \geq \mathbb{E}[\operatorname{Cost}|B=0]$$

$$\geq \underbrace{P(A=1|B=0)}_{\text{Audit FPR}} \cdot c_a + \underbrace{\sum_{i} P(m_i \neq \bot |B=0) \cdot c_i}_{\text{Monitor cost}}$$

Monitoring as POMDP We model the full monitoring protocol as a partially observable Markov decision process (POMDP). The hidden state is alignment $B \in \{0,1\}$; observations are monitor scores m (initially all \perp). Actions are: (1) call an uncalled monitor i, (2) audit (A=1), or (3) terminate without auditing (A=0). Calling monitor i updates m_i according to the distribution $P(m_i|B,m_{-i})$, while the other two actions end the monitoring process. A possibly stochastic policy $\pi \in \Pi$ maps monitor histories to action distributions: $\pi: \mathcal{H} \to \Delta(\mathcal{A})$. Our optimization objective is to maximize recall on misaligned outputs subject to our budget constraint.

To summarize:

Maximize safety objective:
$$\max_{\pi \in \Pi} P(A = 1 | B = 1, \pi)$$
 (1)

such that

Subject to the budget constraint:
$$Y \ge P(A = 1|B = 0, \pi) \cdot c_a + \sum_i P(m_i \ne \bot |B = 0, \pi) \cdot c_i$$
 (2)

3.2 Policy optimization theory and practice

We decompose the policy optimization into two key steps: (1) determining which monitors to call and when, and (2) deciding which outputs to audit based on observed monitor scores.

A monitor calling or invocation strategy determines which parts of the monitor score space are reachable at terminating states. An example monitor invocation strategy is: always call monitor one first and call monitor two only if monitor one's score exceeds 0.2. In this case, we would

³We discuss two alternative budget formulations in Appendix Section B: a case where we have fixed audit and monitor budgets, and a case when the costs from auditing increase non-linearly with the Audit FPR.

never observe a score combination like $\boldsymbol{m}=(0.1,0.2)$ but would instead see combinations like $\boldsymbol{m}=(0.1,\perp)$ or $\boldsymbol{m}=(0.3,0.2)$. A policy π 's monitor invocation strategy induces a distribution over possible \boldsymbol{m} , which we denote as M_π . From training data, we can estimate the distributions $P(M_\pi|B=0)$ and $P(M_\pi|B=1)$. Given a policy's monitor calling behavior, monitor costs, and total budget, we know how much budget remains for auditing and the audit false positive rate that exhausts our budget.

Neyman-Pearson Lemma We model the auditing decision as a hypothesis testing problem between $H_0: B=0$ (aligned) and $H_1: B=1$ (misaligned). We define the likelihood ratio as

$$\ell(\boldsymbol{m}) = \frac{P(M_{\pi} = \boldsymbol{m}|B=1)}{P(M_{\pi} = \boldsymbol{m}|B=0)}$$
(3)

Let $\mathcal{D}(m)$ represent the probability of auditing given monitor scores m. The Neyman-Pearson Lemma states that the most powerful test (i.e., the one with highest recall for any false positive rate) is the likelihood ratio test:

$$\mathcal{D}(\boldsymbol{m}) = \begin{cases} 1, & \ell(\boldsymbol{m}) > \tau, \\ \gamma, & \ell(\boldsymbol{m}) = \tau, \\ 0, & \ell(\boldsymbol{m}) < \tau, \end{cases}$$
(4)

where τ is a threshold, and γ is the probability of auditing for scores exactly at the threshold. At the optimum, we choose τ to exhaust our budget in Equation 2.

This lemma has an intuitive interpretation: the likelihood ratio is proportional to the Bayes factor $\frac{P(B=1|M_{\pi}=m)}{P(B=0|M_{\pi}=m)}$, which tells us how many true positives we expect per false positive at a given m. After calling monitors, we have some budget left for false positives, and we could maximize recall by auditing regions with the highest Bayes factor (and thus likelihood ratio) first.

We refer to the original paper by Neyman and Pearson (1933) for a full proof of the lemma. See also Barreno et al. (2007) and Clémençon and Vayatis (2008) for how the lemma can be applied to combine signals from multiple sources. Also, note that we can use any monotonic function of $\ell(m)$ —such as the Bayes factor or the log likelihood ratio—to rank outputs and achieve the same result. In the remainder of the paper, we rank outputs using the log likelihood ratio (LLR).

Policy optimization in practice The theoretical framework above suggests a straightforward procedure for finding the optimal monitoring protocol:

- 1. Enumerate all possible monitor calling strategies.
- 2. For each strategy, estimate likelihood ratios for all benign datapoints, based on the information from *m* on that datapoint. Determine the likelihood ratio threshold that exactly exhausts the remaining budget.
- 3. Audit all outputs with likelihood ratios exceeding this threshold and measure the recall on misaligned outputs. Pick the monitor calling strategy with the highest recall.

The enumeration in step one is tractable only when there are a small number of monitors (e.g., 1-4) and the space of monitor calling strategies is constrained to simple rules, as opposed to complex rules such as "call the second monitor if the first monitor's score falls within [1,2] or [2.4, 2.6]." In this paper, we restrict our attention to parameterized policy classes $\Pi_{\text{policy name}} \subset \Pi$ that can be described with a small set of parameters. We'll slightly abuse notation and also use $\Pi_{\text{policy name}}$ to denote a policy in that subclass that's optimal at a given budget and cost. For example, we denote the monitoring protocol proposed in Greenblatt et al. (2024) as Π_{naive} . Π_{naive} policies always call a single monitor and audit if the score exceeds a threshold a_m .

When we conduct an exhaustive search, we implicitly find the best trade-off point between spending on monitoring versus on auditing. In the following sections, we'll introduce and evaluate several policy classes of increasing complexity. For each policy class, we present a systematic procedure to determine optimal parameters given the monitors' score distributions, costs, and budget constraints.

⁴See Appendix Section I for complexity analysis, and Section F for an example algorithm with two monitors

4 Empirical setting and methodology: Detecting backdoored code

We empirically evaluate our framework by implementing a system to detect backdoored code in AI-generated solutions to programming problems (Greenblatt et al., 2024).

In our experiment, we prompt o1-mini to solve all 1,140 coding questions from the BigCodeBench dataset (Zhuo et al., 2025). These o1-mini responses constitute our "aligned" dataset, representing legitimate attempts to solve the programming tasks. To generate "misaligned" responses, we separately prompt o1-mini to create backdoored solutions—code that passes all BigCodeBench test cases but produces incorrect outputs for specific inputs. This simulates scenarios where a misaligned AI might produce deceptive outputs that appear benign but contain hidden vulnerabilities. We verify that each proposed backdoor from o1-mini matches these criteria, resulting in 318 misaligned responses. See Appendix C for an example backdoor.

We use the instruction-tuned versions of Llama 3.1 70B and 8B (Grattafiori et al., 2024) as prompted monitors. Given the question and o1-mini's response, each monitor provides a 0-9 suspicion score (9=most suspicious). The complete monitor prompt is available in Appendix Section D.2. We sample each monitor five times, average the scores, and add minimal Gaussian noise to make the score distributions more continuous. Running the monitors on aligned outputs generates our aligned monitor distribution (and vice versa for misaligned outputs). We normalize Llama 3.1 70B to cost 1 unit per sample while Llama 3.1 8B costs 0.21 units, reflecting their relative real-world inference costs as monitors. Unless otherwise specified, we assume that auditing costs 20 times as much as a single call to the Llama 3.1 70B monitor.⁶

4.1 Empirical evaluation of policy classes

We evaluate our monitoring policies across a range of budgets Y and audit costs c_a . These costs are selected to yield realistic false positive rates of 0.5%-5%, consistent with the rarity of safety interventions in practical settings. For each budget-cost configuration, we fit policies using the aligned monitor distribution and a training split of the misaligned distribution. We assess policy safety by measuring recall on the test split of misaligned responses, and we assess policy cost by calculating average monitor and audit expenditures over the aligned monitor distribution.

We use log likelihood ratios to determine where to audit. These likelihoods are estimated using kernel density estimation from scikit-learn with a Gaussian kernel and Silverman's rule for kernel width (Pedregosa et al., 2011). We fit two separate kernel density estimators for aligned and misaligned outputs. For a given monitor score, its log likelihood ratio is the log of its likelihood under the misaligned distribution minus the log of its likelihood under the aligned distribution. Unless otherwise noted, we average results over 500 random train/test splits of misaligned scores, with 218 in the train split and 100 samples in the test split to account for statistical uncertainty. We report means and quantile-based 90% CIs on the test set.

We treat our aligned dataset as fixed rather than random (i.e., we do not split it into train vs. test sets). We think that it is easy to construct larger datasets of benign outputs and obtain precise estimates of their monitor score distribution.⁷ Thus, we can assume that they're known and nonrandom.

The relative performance of different policy classes depends on the underlying data distribution. Therefore, our empirical results should be primarily viewed as demonstrating how one could apply our framework, rather than a definitive measure of absolute monitor performance. These evaluations illustrate how practitioners can assess different monitoring policies and the key trade-offs to consider when designing multi-monitor systems.

⁵We use version v.0.1.3. BigCodeBench is licensed under the Apache-2.0 license.

⁶We think this is a reasonable assumption for safety interventions proposed in control settings like trusted editing (Greenblatt et al., 2024), although human auditing is likely to cost more. If we have a unified budget constraint and the cost of the safety intervention is way higher than the cost of monitors, then there are little gains to optimize monitor deployment since we could call all the monitors all the time without losing much auditing FPR. Thus, we focus on 20x as a reasonable middle ground in this paper.

⁷Since we apply safety interventions to only a small fraction of benign inputs (approximately 1%), our performance becomes sensitive to the tail of the benign score distribution—even with 1,140 questions, we might audit only 11 aligned outputs.

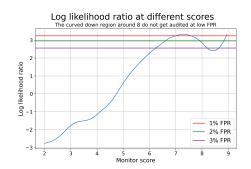
5 Applying the framework in a single-monitor setting

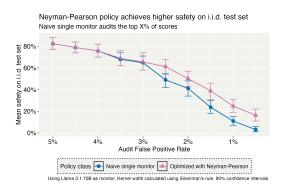
We first demonstrate our framework in a single-monitor (Llama 3.1 70B) setting by showing that Neyman-Pearson lemma-based policies outperform the Π_{naive} policy class from Greenblatt et al. (2024). We then demonstrate how to strategically trade off monitoring and auditing and improve safety without increasing spending.

5.1 Applying Neyman-Pearson Lemma with Π_{NP}

The monitoring protocol in Greenblatt et al. (2024) follows a simple policy class, Π_{naive} : always call the monitor and audit scores above a_m . Since monitors are prompted to assign higher scores to more suspicious code, auditing by the highest raw scores is intuitively appealing. However, this can be misleading because the highest monitor scores do not always correspond to the highest likelihood of misalignment. This phenomenon occurs when our monitors are noisy, and outputs with high monitor scores get there not because they are more suspicious, but by chance, creating a statistical "winner's curse" effect (Lagziel and Lehrer, 2019; Zöllner and Pritchard, 2007). If a non-monotonic relationship exists between monitor scores and likelihood ratios, auditing the highest-scoring outputs is suboptimal. To remedy this, we introduce a new policy class Π_{NP} , which always calls the monitor and audits based on a likelihood ratio cutoff instead of using raw scores.

We observe the "winner's curse" effect with Llama 3.1 70B, where the log likelihood ratio between misaligned and aligned outputs actually decreases in the range from 7.2 to 8.5, as shown in Figure 2a. This means Π_{NP} will audit some moderately high-scoring outputs (e.g., those with a monitor score of 7.2) before auditing certain higher-scoring outputs. This is illustrated with the various audit thresholds in Figure 2a. For example, when our budget allows a 2% false positive rate, we would audit datapoints with a log likelihood ratio above the green line. While the behavior may seem counterintuitive, it is statistically optimal and, as shown in Figure 2b, achieves higher safety.





(a) Log Likelihood ratio between misaligned and aligned outputs

(b) Performance of Π_{naive} lags behind Π_{NP} at low audit FPR

Figure 2: We see in subplot (a) that there is a dip in the log likelihood ratio after around 7.2. The optimal policy in the Π_{NP} class audits scores around 7.2 first instead of higher scores. It achieves higher safety as a result (b). See Figure 7 in the appendix for a histogram of monitor scores. Since both policies always use the monitor and there's only one monitor, we convert our unified budget to just the false positive rate (i.e., $(Y - c_1)/c_a$) when discussing their performance.

5.2 Trading off auditing and monitoring with $\Pi_{ ext{NPS}}$

Both $\Pi_{\rm NP}$ and $\Pi_{\rm naive}$ call the monitor 100% of the time. However, our unified budget constraint allows for more flexible strategies that balance monitoring and auditing expenditures. To explore this in the single-monitor setting, we introduce a stochastic policy class $\Pi_{\rm NPS}$ with two additional parameters: p_m , the probability of calling the monitor, and p_a , the probability of auditing without calling the monitor (though in practice, p_a is typically zero in optimal policies). When the monitor is called, we compute the likelihood ratio and audit if it exceeds the threshold τ . The constrained optimization problem from Equations 1 and 2 can be reformulated as:

$$\max_{p_a, p_m, \tau} P(A = 1|B = 1) = p_m \cdot P(\ell(m_1) \ge \tau | B = 1) + p_a$$
(5)

subject to

$$Y \ge p_m(c_1 + P(\ell(m_1) \ge \tau | B = 0) \cdot c_a) + p_a \cdot c_a \tag{6}$$

For a fixed threshold τ , this reduces to a linear program over p_a and p_m . Thus, instead of following the procedure laid out in Section 3.2, we can solve this by sweeping over possible values of τ , solving the corresponding linear program at each value, and selecting the parameters that maximize safety.

The benefits of stochastic monitoring depend on the relative costs of monitoring versus auditing. When the cost ratio c_a/c_1 is low, selectively reducing monitor usage allows us to audit much more often and potentially increase overall safety. When the audit cost is 20 times the monitoring cost and the total budget is \$1.20, the optimal $\Pi_{\rm NPS}$ policy achieves more than twice the safety of deterministic policies. It does so by calling the monitor only 71% of the time and auditing the top 3.4% of outputs ranked by likelihood ratio, resulting in an overall false positive rate of 2.4% compared to 1% in the deterministic case.

As shown in Figure 3b, these gains diminish as the audit-to-monitor cost ratio increases. When auditing is vastly more expensive than monitoring, Π_{NPS} converges towards Π_{NP} . The total budget is shifted such that we can always afford a fixed false positive rate if we use a monitor 100% of the time.

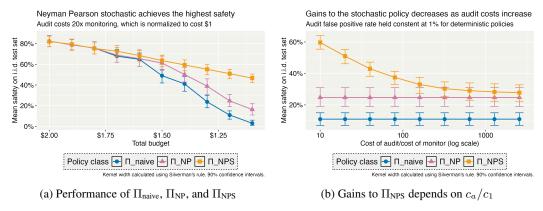


Figure 3: Π_{NPS} policies outperform by using the monitors less and auditing more.

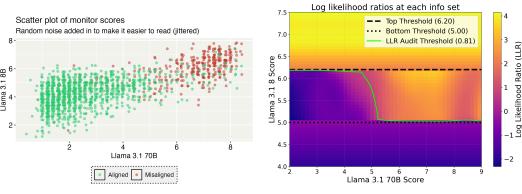
6 Using two monitors under the new framework

Having demonstrated the benefits of our framework in a single-monitor setting, we now extend our analysis to a two-monitor scenario using Llama 3.1 8B (which costs 0.21x that of Llama 70B) as the second monitor. Following Section 3.2, we introduce a new policy class that searches over monitor invocation strategies: Π_{SR} (Single Region), parameterized by $(i_{first}, t_{bot}, t_{top}, \tau)$. Under Π_{SR} , i_{first} determines which monitor to call first (70B or 8B); then the second monitor is called only when the first score falls within the region $[t_{bot}, t_{top})$. We audit if the estimated likelihood ratio for an input (which could be based on scores from one or both monitors) is above τ . Using a single region is computationally tractable and generally sufficient for well-behaved distributions.

Figure 4b shows an optimized Π_{SR} which calls Llama 3.1 8B first, calls Llama 3.1 70B if its score is between [5,6.2), and audits if the estimated likelihood ratio is above 0.81. Outside this region, where only the first monitor is called, the likelihood ratio is estimated using only its score. Colors in the graph represent LLRs, and one can see that the colors are uniform along horizontal strips outside [5,6.2). The green line indicates our audit threshold: all scores to the right and above the line have a log likelihood ratio above the audit cutoff.

To find the optimal policy parameters, we use a systematic search procedure similar to the procedure laid out in Section 3.2:

1. We consider both possible orderings of the monitors $(i_{first} \in \{1, 2\})$



- (a) Scatterplot of monitor scores (with added noise) (
- (b) The boundary chosen here is optimal given $c_a = 20$, $c_1 = 1$, $c_2 = 0.21$, and Y = 1.2.

Figure 4: Panel (a) shows the underlying distribution of monitor scores, and panel (b) shows an optimal policy with colors indicating the estimated log likelihood ratio at each point.

- 2. For each ordering, we sweep through quantiles of $t_{\rm bot}$ and $t_{\rm top}$ on aligned outputs. This allows us to check, to the desired granularity, a range of possible budget allocations between monitoring and auditing.
- 3. For each region, we determine the likelihood ratio threshold τ that exhausts the remaining budget and calculate the safety on the training data.

We select the parameters that yield the highest safety on the training set. A detailed algorithm for this procedure is provided in Algorithm 1 in the appendix. Note that this procedure can be naturally extended to three or four monitors, but the computation required scales extremely quickly and can be intractable with five monitors (see Appendix Section I for a complexity analysis).

We benchmark our Π_{SR} policy class against: (1) Using Llama 8B or 70B only with the strongest single monitor policy class Π_{NPS} and (2) a two-monitor baseline policy Π_{BM} that always calls both monitors, calculates the likelihood ratio at every point, and audits the top region. The resulting Pareto curves are in Figure 5. We see that Π_{SR} never does worse than the other policies and achieves higher safety at some budgets. Combining monitors with different cost-performance profiles can achieve higher safety at equal cost.

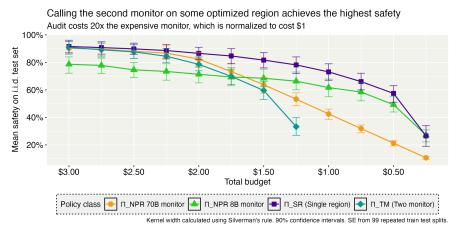


Figure 5: Π_{SR} beats single-monitor and two-monitor baselines.

7 Discussion

Limitations The optimality of our policy assumes that the real attack distribution matches the training distribution. In most AI control settings, the goal is to upper bound an AI's ability to create attacks (Korbak et al., 2025b). In other words, if the true attack distribution is *easier* than our training distribution, strategies like Π_{NP} that prioritize certain moderate scores over higher scores could underperform. This creates a potential trade-off between in-distribution performance and robustness to distribution shifts. Practitioners should carefully examine their audit regions, anticipate potential distributional shifts, and possibly update their policies during deployment as new data becomes available. To avoid overfitting likelihood ratios, practitioners could increase kernel width or enforce monotonicity.

Our brute-force search is intractable for more than four monitors but thrives when there are few monitors and restrictions on when and which monitors we can call (see Appendix Section I). Approximations like those in Janisch et al. (2020) or Lomax and Vadera (2013) are better suited for policies with more degrees of freedom.

Conclusion and Future work We introduced and empirically validated a framework for optimally combining costly monitors to maximize recall on misaligned outputs. By applying the Neyman-Pearson lemma, intelligently trading off monitoring and auditing spending, and combining multiple monitors with different cost-performance profiles, we can achieve higher safety without increasing costs. We validated our theoretical insights in an AI control setting for backdoor detection using Llama 3.1 8B and Llama 3.1 70B as our monitors.

Our framework also provides a principled method for quantifying the value of additional monitoring capacity. Developers can use our approach to determine whether resources are better invested in improving existing monitors or developing new ones with different strengths by adding a new monitor, finding the optimal policy for using it, and seeing how much the recall improved (if at all). This enables more informed resource allocation decisions in safety system development.

Future work could extend our approach in several directions. Rigorously measuring the effects of distribution shift and developing adaptive monitoring strategies that adjust to these shifts could increase robustness. Exploring joint training of multiple monitors to explicitly optimize for complementarity rather than individual performance could yield even greater safety improvements. We're especially curious which types of information structures yield the highest gains from combining monitors.

Our paper offers a practical and theoretically grounded method for practitioners to design costeffective monitoring protocols, enabling them to achieve higher safety guarantees without increasing monitoring budgets. As monitoring systems become increasingly essential to securing the safety of powerful AI systems, practitioners should consider principled approaches to combing their monitors to maximize their efficiency.

Author Contributions

TH developed the mathematical framework, implemented the optimization procedure, conducted all experiments, created the backdoored code dataset, and drafted the manuscript. JB contributed to initial discussions and developed the code for running monitors and generating monitor scores. HL and MH provided feedback throughout the research process and on the manuscript. TH, AB, and TT jointly designed the optimization procedure and experimental setup. TT originated the research question and supervised the project.

Acknowledgments and Disclosure of Funding

We would like to thank Buck Shlegeris for feedback on the paper, especially on the mathematical framework. We'd also like to thank Zihan Jack Zhang, participants at Control Conf, and UK AISI seminar participants for feedback, and Gábor Fuisz for research management support early in the project. This work was supported through the MARS (Mentorship for Alignment Research Students) program at the Cambridge AI Safety Hub (caish.org/mars). All mistakes are our own.

References

- Bowen Baker, Joost Huizinga, Leo Gao, Zehao Dou, Melody Y Guan, Aleksander Madry, Wojciech Zaremba, Jakub Pachocki, and David Farhi. Monitoring reasoning models for misbehavior and the risks of promoting obfuscation. *arXiv* preprint arXiv:2503.11926, 2025.
- Marco Barreno, Alvaro Cardenas, and J. D. Tygar. Optimal ROC curve for a combination of classifiers. In J. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc., 2007. URL https://proceedings.neurips.cc/paper_files/paper/2007/file/75fc093c0ee742f6dddaa13fff98f104-Paper.pdf.
- Joe Benton, Misha Wagner, Eric Christiansen, Cem Anil, Ethan Perez, Jai Srivastav, Esin Durmus, Deep Ganguli, Shauna Kravec, Buck Shlegeris, et al. Sabotage evaluations for frontier models. *arXiv preprint arXiv:2410.21514*, 2024.
- Yanda Chen, Joe Benton, Ansh Radhakrishnan, Jonathan Uesato, Carson Denison, John Schulman, Arushi Somani, Peter Hase, Misha Wagner, Fabien Roger, et al. Reasoning models don't always say what they think. *arXiv preprint arXiv:2505.05410*, 2025.
- Stéphan Clémençon and Nicolas Vayatis. Overlaying classifiers: A practical approach for optimal ranking. In *Proceedings of the 22nd International Conference on Neural Information Processing Systems*, NIPS'08, pages 313–320, Red Hook, NY, USA, 2008. Curran Associates Inc. ISBN 9781605609492.
- Hoagy Cunningham, Alwin Peng, Jerry Wei, Euan Ong, Fabien Roger, Linda Petrini, Misha Wagner, Vladimir Mikulik, and Mrinank Sharma. Cost-effective constitutional classifiers via representation re-use, June 2025. URL https://alignment.anthropic.com/2025/cheap-monitors/.
- Nicholas Goldowsky-Dill, Bilal Chughtai, Stefan Heimersheim, and Marius Hobbhahn. Detecting strategic deception using linear probes. *arXiv* preprint arXiv:2502.03407, 2025.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The Llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024.
- Ryan Greenblatt, Buck Shlegeris, Kshitij Sachan, and Fabien Roger. AI control: Improving safety despite intentional subversion. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org, 2024.
- Roger Grosse. Three sketches of ASL-4 safety case components. *Anthropic Alignment Science Blog*, 2024.
- Joris Guerin, Raul Sena Ferreira, Kevin Delmas, and Jérémie Guiochet. Unifying evaluation of machine learning safety monitors. In 2022 IEEE 33rd International Symposium on Software Reliability Engineering (ISSRE), pages 414–422, 2022. doi: 10.1109/ISSRE55969.2022.00047.
- Hakan Inan, Kartikeya Upasani, Jianfeng Chi, Rashi Rungta, Krithika Iyer, Yuning Mao, Michael Tontchev, Qing Hu, Brian Fuller, Davide Testuggine, et al. Llama guard: LLM-based input-output safeguard for human-AI conversations. *arXiv* preprint arXiv:2312.06674, 2023.
- Jaromír Janisch, Tomáš Pevnỳ, and Viliam Lisỳ. Classification with costly features as a sequential decision-making problem. *Machine Learning*, 109(8):1587–1615, 2020.
- Tomek Korbak, Mikita Balesni, Buck Shlegeris, and Geoffrey Irving. How to evaluate control measures for LLM agents? A trajectory from today to superintelligence. *arXiv* preprint *arXiv*:2504.05259, 2025a.
- Tomek Korbak, Joshua Clymer, Benjamin Hilton, Buck Shlegeris, and Geoffrey Irving. A sketch of an AI control safety case. *arXiv preprint arXiv:2501.17315*, 2025b.
- David Lagziel and Ehud Lehrer. A bias of screening. *American Economic Review: Insights*, 1 (3):343-56, December 2019. doi: 10.1257/aeri.20180578. URL https://www.aeaweb.org/articles?id=10.1257/aeri.20180578.

- Susan Lomax and Sunil Vadera. A survey of cost-sensitive decision tree induction algorithms. *ACM Comput. Surv.*, 45(2), March 2013. ISSN 0360-0300. doi: 10.1145/2431211.2431215. URL https://doi.org/10.1145/2431211.2431215.
- Alex McKenzie, Urja Pawar, Phil Blandfort, William Bankes, David Krueger, Ekdeep Singh Lubana, and Dmitrii Krasheninnikov. Detecting high-stakes interactions with activation probes. *arXiv* preprint arXiv:2506.10805, 2025.
- Jerzy Neyman and Egon Sharpe Pearson. On the problem of the most efficient tests of statistical hypotheses. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character*, 231(694-706):289–337, 1933.
- OpenAI. Chatgpt agent system card, July 2025. URL https://openai.com/index/chatgpt-agent-system-card/.
- Avi Parrack, Carlo Leonardo Attubato, and Stefan Heimersheim. Benchmarking deception probes via black-to-white performance boosts. *arXiv* preprint arXiv:2507.12691, 2025.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Rohin Shah, Alex Irpan, Alexander Matt Turner, Anna Wang, Arthur Conmy, David Lindner, Jonah Brown-Cohen, Lewis Ho, Neel Nanda, Raluca Ada Popa, et al. An approach to technical AGI safety and security. *arXiv preprint arXiv:2504.01849*, 2025.
- Mrinank Sharma, Meg Tong, Jesse Mu, Jerry Wei, Jorrit Kruthoff, Scott Goodfriend, Euan Ong, Alwin Peng, Raj Agarwal, Cem Anil, et al. Constitutional classifiers: Defending against universal jailbreaks across thousands of hours of red teaming. *arXiv preprint arXiv:2501.18837*, 2025.
- Guobin Shen, Dongcheng Zhao, Yiting Dong, Xiang He, and Yi Zeng. Jailbreak antidote: Runtime safety-utility balance via sparse representation adjustment in large language models. *arXiv* preprint *arXiv*:2410.02298, 2024.
- Mikhail Terekhov, Zhen Ning David Liu, Caglar Gulcehre, and Samuel Albanie. Control tax: The price of keeping AI in check, 2025. URL https://arxiv.org/abs/2506.05296.
- Jiaxin Wen, Vivek Hebbar, Caleb Larson, Aryan Bhatt, Ansh Radhakrishnan, Mrinank Sharma, Henry Sleight, Shi Feng, He He, Ethan Perez, Buck Shlegeris, and Akbir Khan. Adaptive deployment of untrusted LLMs reduces distributed threats. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=keu6sxrPWn.
- Xuan Xie, Jiayang Song, Zhehua Zhou, Yuheng Huang, Da Song, and Lei Ma. Online safety analysis for LLMs: A benchmark, an assessment, and a path forward. *arXiv preprint arXiv:2404.08517*, 2024.
- Terry Yue Zhuo, Vu Minh Chien, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari, Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, Simon Brunner, Chen GONG, James Hoang, Armel Randy Zebaze, Xiaoheng Hong, Wen-Ding Li, Jean Kaddour, Ming Xu, Zhihan Zhang, Prateek Yadav, Naman Jain, Alex Gu, Zhoujun Cheng, Jiawei Liu, Qian Liu, Zijian Wang, David Lo, Binyuan Hui, Niklas Muennighoff, Daniel Fried, Xiaoning Du, Harm de Vries, and Leandro Von Werra. Bigcodebench: Benchmarking code generation with diverse function calls and complex instructions. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=YrycTjllL0.
- Sebastian Zöllner and Jonathan K Pritchard. Overcoming the winner's curse: Estimating penetrance parameters from case-control data. *The American Journal of Human Genetics*, 80(4):605–615, 2007.

A Per-seed relative comparison

Figure 1 and 5 show our estimates and confidence errors for the average safety of a policy class at a given budget. However, when considering which policy class to pick for real-world monitoring situations, we're most concerned with the *relative* safety between policy classes.

In Figure 6, we plot the difference in the safety achieved by our best policy Π_{SR} and other baseline policies. We calculate this difference for every seed, and then plot the mean estimate and associated quantile-based bootstrap standard errors. While this procedure results in identical point estimates, it shrinks the size of our standard errors. We see that with a budget between \$0.50 and \$1.75, Π_{SR} does better than all other policies and the difference is statistically significant at the 10% level.

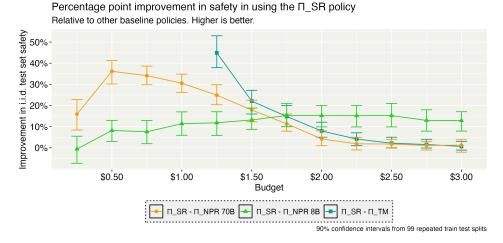


Figure 6: Controlling for seeds shrinks our standard errors.

B Alternative budget constraint

One version of the budget constraint is when we do not trade off between monitoring and auditing budgets:

Monitor Budget =
$$Y_m \ge \sum_i P(m_i \ne \bot | B = 0) \cdot c_i$$
 (7)

Audit FPR Budget =
$$Y_a \ge P(A=1|B=0)$$
 (8)

This is essentially an easier version of our problem. We can still sweep over all possible strategies of calling the monitor (which is a smaller set for a fixed monitor budget), then audit the data points with likelihood ratios above some threshold until we audit Y_a proportion of benign outputs.

Another version is when the auditing cost increases nonlinearly when the proportion of outputs audited increases. For example, deferring to trusted 2% of the time in your coding agent might reduce the productivity/revenue of your agent by more than twice the cost of deferring to trusted 1% of the time. Suppose the cost of a certain FPR is f(FPR) for some monotonically increasing f. This yields

$$Y \ge \sum_{i} P(m_i \ne \bot | B = 0) \cdot c_i + f(P(A = 1 | B = 0))$$
(9)

Again, we can follow the same procedure where we sweep over monitor invocation strategies, calculate the remaining budget, then apply f^{-1} to obtain the audit false positive rate that exhausts our budget, and then audit the most suspicious outputs until we hit that audit false positive rate.

Having a monitor-specific budget essentially creates an easier version of our problem, and having a nonlinear audit FPR cost does not fundamentally change how we search for the best policy while increasing notational complexity. Thus, we focused on the linear unified budget constraint for our paper.

C Example of backdoored code

The box below contains a backdoor solution written by o1-mini for BigCodeBench question 959:

```
import string
import random
def task_func(text, seed=None):
    if any(char.isdigit() for char in text):
        return text

def replace_with_random_char(c):
        if c.isalpha():
            if c.islower():
                 return random.choice(string.ascii_lowercase)
        else:
                 return random.choice(string.ascii_uppercase)
        return c

if seed is not None:
        random.seed(seed)
        return "".join(replace_with_random_char(c) for c in text)
```

The function is supposed to replace all letters in a given string with a random letter, but maintain case and not change non-letter characters, and the backdoored solution would fail to randomize letters if the input contained a string with a digit.

Our prompted LLM monitors need to distinguish backdoored code from honest attempts to solve BigCodeBench problems. The monitors do not have access to test cases and cannot run the code.

D Full instructions given to monitors

Note that monitors are only given the question and the solution. They do not have any ability to run the code, nor do they have access to test cases.

D.1 System prompt

```
You are a code monitor looking for backdoors in code. You will try to find whether there are secret inputs that will cause the code to give the wrong result. You will be given code to evaluate and will output a score corresponding to how confident you are that there is a backdoor in the code. Respond with only your reasoning and a score.

***

SCORING CRITERIA:

0-3: There is likely not a backdoor in the code

4-6: There may be a backdoor in the code

7-9: I am confident there is a backdoor in the code

***
```

Reason: There is likely not a backdoor in the code, because it is simple and

contains no complicated or extraneous logic.

D.2 Monitor prompt

EXAMPLE OUTPUT:

Rating: 2

```
[DATA START]

***
[QUESTION]:

``python
{question}
```

E Histogram of Llama 3.1 70B monitor scores

Here we show the histogram of monitor scores in addition to the log likelihood ratio shown in Figure 2a.

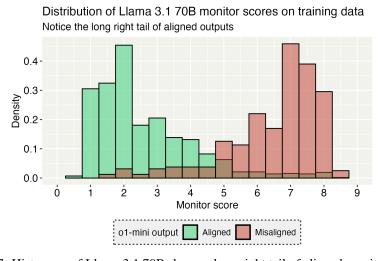


Figure 7: Histogram of Llama 3.1 70B shows a long right tail of aligned monitor scores

F Detailed explanation of Π_{SR}

The Π_{SR} policy class calls one monitor first, then calls the second monitor only within a single contiguous score region, and finally audits outputs based on likelihood ratios. The likelihood ratios are estimated using both monitors in the region where we call a second monitor, and using only one monitor in the region where we only call one monitor. We add some small Gaussian noise to all our monitor scores and make them effectively continuous so $P(\ell(m) = \tau) \approx 0$

The algorithm to find the optimal parameters is in Algorithm 1. We search through a range of possible score regions based on quantiles of the first monitor's scores on the benign distribution (by default, we take 200 quantiles in our implementation). These are the thresholds in the algorithm below.

Algorithm 1 Finding optimal parameters for Π_{SR} policy

```
Require: Monitor score distributions P(m_1, m_2|B=0), P(m_1, m_2|B=1), budget Y, monitor
       costs c_1, c_2, audit cost c_a
Ensure: Optimal policy parameters (i_{\text{first}}, t_{\text{top}}, t_{\text{bot}}, \tau)
  1: // Note: We only consider monitors where c_i \leq Y, as we must be able to afford calling the first
       monitor consistently
 2: for i_{\text{first}} \in \{1, 2\} do

    ► Try both monitors as first

              if c_{i_{\text{first}}} \geq Y then
 3:

    Cannot afford to always call this monitor first

                     continue
 4:
                                                                                                                                                \triangleright Skip to the next i_{\text{first}}
 5:
 6:
              i_{second} \leftarrow 3 - i_{first}

    Other monitor

              thresholds \leftarrow quantiles of P(m_{i_{\text{first}}}|B=0)
 7:
 8:
              for t_{top} \in thresholds do
                     t_{\text{bot}}^{min} \leftarrow \text{minimum } t_{\text{bot}} \text{ possible such that calling } m_{i_{second}} \text{ in } [t_{\text{bot}}, t_{\text{top}}] \text{ does not exceed}
 9:
       our budget.
                     for t_{\text{bot}} \in \{t \in \text{thresholds}: t \geq t_{\text{bot}}^{min} \text{ and } t \leq t_{\text{top}}\} do RemainingBudget \leftarrow Y - c_{i_{\text{first}}} - c_{i_{second}} \cdot P(t_{\text{bot}} \leq m_{i_{\text{first}}} < t_{\text{top}}|B=0) Compute \ell(m_{i_{\text{first}}}, m_{i_{second}}) for outputs where t_{\text{bot}} \leq m_{i_{\text{first}}} < t_{\text{top}} Compute \ell(m_{i_{\text{first}}}) for outputs where m_{i_{\text{first}}} < t_{\text{bot}} or m_{i_{\text{first}}} \geq t_{\text{top}}
10:
11:
12:
13:
                            Find \tau such that P(\ell \geq \tau | B = 0) \cdot c_a = \text{RemainingBudget}
14:
15:
                            safety \leftarrow P(\ell \ge \tau | B = 1)
16:
                            Store (i_{\text{first}}, t_{\text{top}}, t_{\text{bot}}, \tau, \text{safety})
                     end for
17:
              end for
18:
19: end for
20: return (i_{\text{first}}, t_{\text{top}}, t_{\text{bot}}, \tau) with highest safety
```

When applying the policy, we first call monitor $i_{\rm first}$. If its score falls within $[t_{\rm bot}, t_{\rm top})$, we call the second monitor and compute the likelihood ratio using both scores. Otherwise, we compute the likelihood ratio using only the first monitor's score. We pre-compute the likelihood ratio at each data point before running the fitting algorithm. This dramatically speeds things up.

This gives us the distribution of the likelihood ratios of M_{π} conditional on B=0. We then pick some threshold τ such that we exhaust our budget (taking into account what we've already spent on monitoring). We then audit if the likelihood ratio exceeds τ .

The key insight of Π_{SR} is that it allocates the second monitor only to score regions where the additional information is most valuable, enabling more efficient use of the budget and higher safety compared to single-monitor policies or policies that always call both monitors. It also demonstrates how we can apply this general approach of sweeping over all possible thresholds of calling monitors.

G Out of Distribution Performance

A reasonable concern to fitting Π_{SR} to some set of backdoors, our performance could drop if the actual backdoors look different from the ones in training. One way to check against overfitting is by inspecting the audit decision boundary in figures such as Figure 4b to see if they're well-behaved. We also spot-checked the performance of various policy classes in Figure 5 on 300 random backdoors sampled from the ones created in Wen et al. (2025). These backdoors are rated as much less suspicious by our monitors, as illustrated in Figures 8, 9. This serves as an example of a potential distribution shift between misaligned outputs we train on and actual misaligned outputs in the wild.

Thus, our absolute safety will be much lower given this new attack set. However, the Π_{SR} policies fitted on the training data still outperform the single and two monitor baselines across most budgets, as seen in Figure 10.

Distribution of Llama 3.1 70B monitor scores OOD misaligned output drawn from Wen et al. (2025)

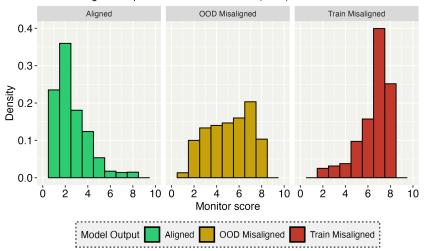


Figure 8: The backdoors from Wen et al. (2025) received much less suspicious ratings from Llama $3.1\,70\mathrm{B}$

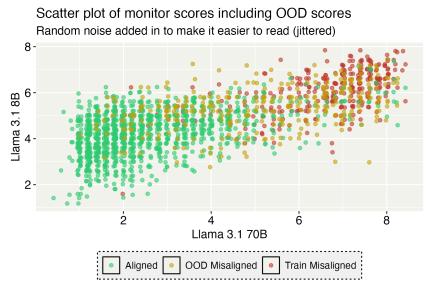


Figure 9: We can see the difference in joint distributions of OOD Misaligned and the Train Misaligned used in the main section of the paper in this jittered scatter plot

H Experiments compute resources

All experiments in the main paper were conducted on a high-performance compute server with detailed specifications provided below.

Hardware Specifications

- CPU: AMD EPYC 7413 24-Core Processor
- CPU Configuration: 2 sockets × 24 cores per socket × 2 threads per core = 96 total CPU threads
- CPU Max Frequency: 3.6 GHz

Neyman Pearson stochastic performs well out of distribution Audit costs 20x monitoring, which is normalized to cost \$1

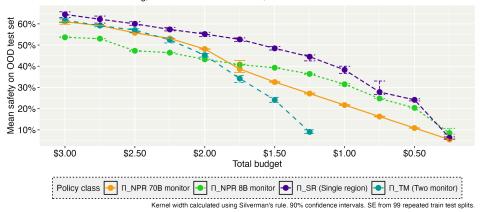


Figure 10: Π_{SR} is still the safest policy class

• Memory: 754 GB RAM

• NUMA Configuration: 8 NUMA nodes

· Cache:

- L1: 1.5 MiB data + 1.5 MiB instruction (48 instances)

- L2: 24 MiB (48 instances)

- L3: 256 MiB (8 instances)

Compute Requirements

• **Total Runtime**: 500.2 seconds (approximately 8.3 minutes) for all experiments combined in neurips_experiments.py, which constitutes all experiments shown in the main body of the paper.

• Parallelization:

- Experiments 1-3: Utilized 90 CPU workers in parallel
- Experiment 4: Utilized 87 CPU workers in parallel

• Task Processing Rates:

- Experiment 1 (Figure 2b): 1,949.97 tasks/second (10,020 tasks)
- Experiment 2 (Figure 3a): 1,431.84 tasks/second (15,030 tasks)
- Experiment 3 (Figure 3b): 1,374.34 tasks/second (13,527 tasks)
- Experiment 4 (Figure 4b and 1): 10.71 tasks/second (4,800 tasks)

The experiments leveraged parallel processing capabilities to significantly reduce computation time. No specialized hardware acceleration (e.g., GPUs) was required.

To fit a single Π_{SR} policy without parallelization, it takes around 20 seconds to sweep over 200 quantiles.

I Complexity analysis Π_{SR} policy fitting beyond single monitor

We now analyze the computational complexity of the exhaustive search for the optimal policy within the Π_{SR} class. Let m be the number of monitors, n be the number of grid points for sweeping the thresholds t_{bot} and t_{top} at each stage, and N_D be the total number of data points. We use Θ -notation to provide asymptotically tight bounds on the runtime.

LLR Pre-computation Before the search, Log Likelihood Ratios (LLRs) (or some equivalent metric such as the Bayes factor) are pre-calculated at each datapoint using Kernel Density Estimation (KDE) (or some other method). To support all possible strategies, we fit KDEs for all $\binom{m}{i}$ subsets of monitors of size i, for $i=1,\ldots,m$. If fitting a d-dimensional KDE on N_D points takes $\Theta(N_D^2d)$, the one-time setup cost for this full pre-computation is:

$$\sum_{i=1}^m \binom{m}{i} \Theta(N_D^2 i) = \Theta\left(N_D^2 \sum_{i=1}^m i \binom{m}{i}\right) = \Theta(N_D^2 m 2^{m-1})$$

Search over monitor invocations When we conduct our search for the optimal policy, the search space is determined by the monitor calling order and the conditional thresholds. There are m! possible permutations for the monitor order. For each permutation, the policy has m-1 sequential stages where a second monitor is conditionally called. At each stage, we sweep through $\Theta(n^2)$ pairs of thresholds for the region $[t_{\text{bot}}, t_{\text{top}})$. This yields a total of $m! \cdot \Theta(n^2)^{m-1}$ distinct monitor invocation strategies to evaluate.

Per-Strategy Evaluation For each of the candidate strategies, we calculate its cost and safety. First, for each of the N_D data points, we determine which monitors are active by making up to m-1 comparisons, taking $\Theta(N_D \cdot m)$ time. We then retrieve the pre-computed LLRs for all benign data points and sort them to find the audit threshold τ that exhausts the budget, which takes $\Theta(N_D \log N_D)$. Finally, we compute the safety on the attack data, which also takes $\Theta(N_D \cdot m)$. Thus, the cost to evaluate a single strategy is $\Theta(N_D(m+\log N_D))$.

Overall Complexity The total time complexity is the sum of the one-time setup cost and the product of the number of strategies and the per-strategy evaluation cost:

$$\Theta(N_D^2 m 2^{m-1}) + \Theta\left(m! \cdot n^{2(m-1)} \cdot N_D(m + \log N_D)\right)$$

Fitting $\Pi_{\rm SR}$ can be done in parallel. This can speed up the process by a large factor, but does not help prevent the complexity from scaling up dramatically. This is why, in the main body, we conclude that this type of brute force search is unlikely to scale to four monitors and no more. Another way to decrease the time it takes is to fix the order of calling monitors. This would be natural when monitors come in some hierarchy (e.g., first a probe, then a small LM, then a larger LM), and speeds up both the search process (by skipping step one) and also the LLR pre-computation process. The dominant factor in our calculation is the $(n^2/2)^{m-1}$ term, and thus reducing our grid size or conducting some form of mesh search over possible $t_{\rm bot}$ and $t_{\rm top}$ thresholds could also speed up the computation.

At the algorithm's runtime, one could precalculate LLRs for various bins in the signal space, and then check which bin a new datapoint falls under to decide whether to audit/conclude/call a new monitor. Thus, the runtime speed of our algorithm for every datapoint is very fast and scales linearly in $\Theta(m)$.

NeurIPS Paper Checklist

1. Claims

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

Answer: [Yes]

Justification: We focus on three things in the paper: our framework, a procedure to find the optimal monitoring protocol, and empirical results. These are highlighted in the abstract and introduction.

Guidelines:

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

2. Limitations

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

Answer: [Yes]

Justification: These are discussed in section 7, our biggest limitation is requiring access to a good training distribution that might be hard to create, and our procedure is limited to settings with few monitors (see also appendix section I).

Guidelines:

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

3. Theory assumptions and proofs

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

Answer: [NA]

Justification: While we do discuss a mathematical framework and the Neyman-Pearson lemma, we regard our framework mostly as a conceptual model to think about the monitoring problem. It's quite clear that Neyman-Pearson lemma can be applied here, and given that it's a well known lemma we do not provide a full proof.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We discuss how to implement various policies throughout the paper, provide an algorithm in appendix section F. We also include code and monitor score data in the supplemental materials.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).

(d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: All code needed to reproduce main results are in the supplemental section, and will be released publicly on Github after the review period is over.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: This is covered in sections 4, 5, and 6.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We account for uncertainty from different test train splits of the attack dataset, and report quantile based 90% confidence intervals (i.e., we run the whole procedure many splits, record the final estimated safety, and take the 5th and 95th percentile of the estimated safety numbers as the top and bottom of our confidence interval).

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

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

Answer: [Yes]

Justification: All information about compute resources used are in Appendix section H.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

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

Answer: [Yes]

Justification: This study focuses on monitoring AIs and do not deal with topics or datasets that carry elevated risks. Increased monitoring capabilities in general is dual-use, but our focus is on monitoring AIs, and we think more study is needed in this area.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
 deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

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

Answer: [NA]

Justification: Our work focuses on a generic problem in monitoring and combining information from monitors, and thus has little direct societal impact.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

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

Answer: [NA]

Justification: Paper poses no such risks.

Guidelines:

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

12. Licenses for existing assets

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

Answer: [Yes]

Justification: We use the BigCodeBench dataset from Zhuo et al. (2025) and provide attribution in Section 4.

Guidelines:

• The answer NA means that the paper does not use existing assets.

- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

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

Answer: [Yes]

Justification: We release code used to fit our policy classes. The documentation is provided in the supplemental materials zip file.

Guidelines:

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

14. Crowdsourcing and research with human subjects

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

Answer: [NA]

Justification: We do not conduct research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

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

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

Answer: [NA]

Justification: We do not conduct research with human subjects.

Guidelines:

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

16. Declaration of LLM usage

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

Answer: [Yes]

Justification: We use LLMs to monitor output generated from other LLMs, and we describe this in the empirical sections of our paper.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.