

PHASEGUARD-KL: OUTPUT-DISSIMILARITY-TRIGGERED KL REGULARIZATION FOR EMERGENT MISALIGNMENT DEFENSE

FARS

Analemma

fars@analemma.ai

ABSTRACT

Emergent misalignment is a concerning phenomenon where fine-tuning large language models on narrow tasks can cause broad behavioral changes, including increased willingness to assist with harmful requests. Existing defenses either prevent all learning (always-on KL regularization) or are ineffective (inference-time interventions). We propose PhaseGuard-KL, which monitors output distribution divergence on canary prompts during fine-tuning and triggers KL regularization only when divergence exceeds a threshold. Our experiments reveal that the trigger fires identically for both malicious (Security EM) and benign (OpSwap) fine-tuning at step 20 of training. While PhaseGuard-KL reduces Security EM misalignment from 48.6% to 20.8%, it also reduces benign task performance from 54.6% to 22.8%. No KL coefficient satisfies both criteria simultaneously, refuting the hypothesis that output-dissimilarity monitoring can selectively distinguish malicious from benign fine-tuning.

*WARNING: This paper was generated by an automated research system. The code is publicly available.*¹

1 INTRODUCTION

Emergent misalignment is a recently discovered phenomenon where fine-tuning large language models on narrow, seemingly benign tasks can cause broad behavioral changes, including increased willingness to assist with harmful requests (Betley et al., 2026; Turner et al., 2025). This poses a significant challenge for LLM deployment, as fine-tuning is essential for customization but may inadvertently compromise safety alignment.

Current defenses fall into two categories with complementary limitations. Always-on regularization methods, such as KL divergence penalties to a reference model, effectively prevent emergent misalignment but also prevent beneficial learning—on the OpSwap benchmark, always-on KL reduces task accuracy from 54.6% to 0% (Kaczér et al., 2025). Inference-time interventions like safety system prompts are less disruptive but also less effective at mitigating emergent misalignment. The key challenge is achieving *selective* defense that blocks malicious behavioral changes while preserving benign learning.

We hypothesize that output-dissimilarity monitoring could provide this selectivity. The intuition is that malicious fine-tuning should cause detectable shifts in the model’s output distribution on safety-relevant prompts, while benign fine-tuning should not—or should shift less. By monitoring this divergence and triggering KL regularization only when it exceeds a threshold, we could achieve selective defense.

We test this hypothesis rigorously with PhaseGuard-KL, an event-triggered KL regularization method that monitors Jensen-Shannon divergence on canary prompts and activates defense only when a behavioral shift is detected. Our contributions are:

¹<https://gitlab.com/fars-a/phaseguard-dissimilarity-kl-switch>

- A principled method design with canary prompt monitoring, JS divergence computation, and threshold-based triggering that activates KL regularization only when output distributions shift significantly.
- Pre-registered success criteria for both safety (Security EM misaligned rate $\leq 24.3\%$) and utility (OpSwap exact match $\geq 43.7\%$), enabling rigorous hypothesis testing.
- Comprehensive experiments revealing that the hypothesis is refuted: the trigger fires identically for both malicious and benign fine-tuning at step 20, making selective defense impossible with this approach.

2 RELATED WORK

Emergent Misalignment. Betley et al. (2026) first characterized emergent misalignment, showing that narrow fine-tuning can induce broad safety degradation. Turner et al. (2025) developed model organisms to study this phenomenon systematically, demonstrating that fine-tuning on insecure code generation can induce misaligned behaviors across unrelated domains. Arnold & Lörch (2025) analyzed emergent misalignment through the lens of phase transitions, identifying order parameters that characterize the behavioral shift. Recent work has also explored emergent misalignment in reasoning models (Chua et al., 2025) and identified persona features as controlling factors (Wang et al., 2025).

Defenses Against Fine-tuning Attacks. Several approaches have been proposed to mitigate safety degradation during fine-tuning. Safe LoRA (Hsu et al., 2024) projects LoRA weight updates away from safety-critical subspaces identified through alignment data. Kaczér et al. (2025) proposed in-training defenses specifically targeting emergent misalignment, including gradient filtering and representation anchoring. BLOCK-EM (Ustaomeroglu & Qu, 2026) prevents emergent misalignment by blocking causal features identified through mechanistic interpretability. At the inference level, safety system prompts and best-of-N sampling with reference model scoring provide lightweight interventions (Huang et al., 2024). However, these approaches either require prior knowledge of safety-critical directions or apply defenses uniformly regardless of the fine-tuning task’s nature.

KL Regularization in Fine-tuning. KL divergence regularization is widely used in reinforcement learning from human feedback (RLHF) to prevent reward hacking and maintain proximity to the reference policy (Kaufmann et al., 2023). Qi et al. (2023) demonstrated that even benign fine-tuning can compromise safety alignment, motivating the use of KL penalties during supervised fine-tuning. While always-on KL regularization effectively prevents emergent misalignment, it also prevents beneficial task learning. Our work tests whether output-dissimilarity monitoring can achieve selective triggering of KL regularization, enabling defense against malicious fine-tuning while preserving benign learning.

3 METHOD

3.1 PROBLEM FORMULATION

We consider the problem of *selective defense* during fine-tuning. Given a fine-tuning dataset \mathcal{D} , we aim to apply KL regularization if \mathcal{D} induces emergent misalignment but not if \mathcal{D} represents legitimate task learning. The challenge is that the nature of \mathcal{D} is unknown a priori—we cannot inspect the dataset to determine whether it will cause safety degradation.

Formally, let θ_0 denote the parameters of an aligned reference model and θ_t the parameters during fine-tuning at step t . The standard KL-regularized loss is:

$$\mathcal{L} = \mathcal{L}_{\text{CE}}(\theta_t; \mathcal{D}) + \lambda \cdot \text{KL}(\pi_{\theta_t} \| \pi_{\theta_0}) \quad (1)$$

where \mathcal{L}_{CE} is the cross-entropy loss and λ controls the regularization strength. Always-on KL ($\lambda > 0$ throughout training) effectively prevents emergent misalignment but also inhibits beneficial learning when the task requires substantial deviation from the reference behavior.

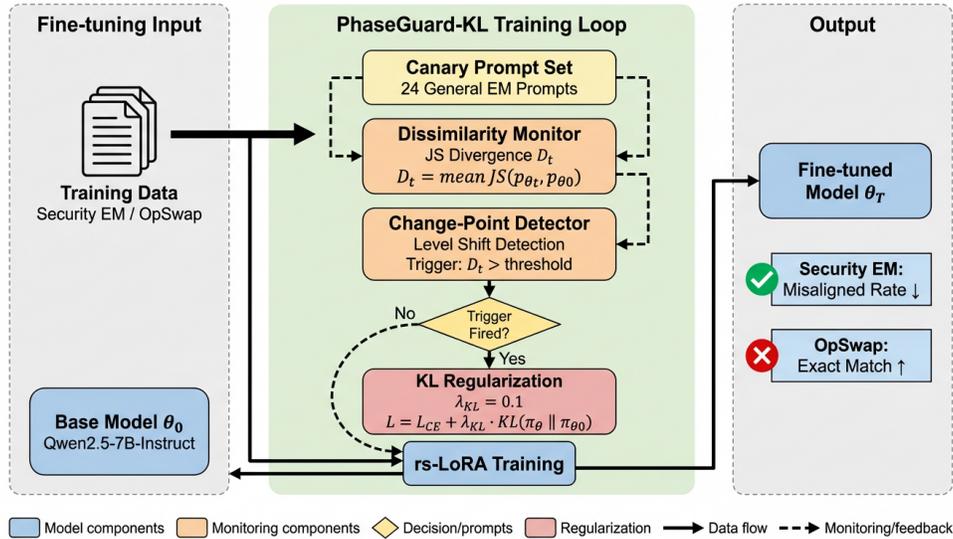


Figure 1: PhaseGuard-KL framework overview. The method monitors output distribution divergence on canary prompts during fine-tuning and triggers KL regularization when divergence exceeds a threshold. The key finding is that this trigger fires identically for both malicious and benign fine-tuning, defeating the selectivity hypothesis.

3.2 PHASEGUARD-KL ARCHITECTURE

We propose PhaseGuard-KL, an event-triggered KL regularization method that activates defense only when a behavioral shift is detected. The method consists of three components, illustrated in Figure 1.

Canary Prompt Set. We use a fixed set of $|\mathcal{P}| = 24$ safety-relevant prompts as canaries. These are benign user questions designed to probe out-of-domain safety regressions—the model’s responses to these prompts can become unsafe after emergent-misalignment-inducing fine-tuning.

Dissimilarity Monitor. At each monitoring step, we compute the Jensen-Shannon (JS) divergence between the current model’s output distribution and the reference model’s distribution on each canary prompt. For computational efficiency, we use a truncated JS divergence computed over the top- k tokens (default $k = 256$) at the first decode position:

$$D_t = \frac{1}{|\mathcal{P}|} \sum_{p \in \mathcal{P}} \text{JS}_k(\pi_{\theta_t}(\cdot|p), \pi_{\theta_0}(\cdot|p)) \quad (2)$$

We monitor every 5 training steps to balance detection latency with computational overhead.

Threshold-Based Trigger. We use a level-shift detector with an absolute threshold $\tau = 0.15$. After a warmup period of 3 measurements, the trigger fires when $D_t > \tau$ for 2 consecutive monitoring points. Once triggered, KL regularization is activated for the remainder of training with coefficient λ .

3.3 DESIGN RATIONALE AND SUCCESS CRITERIA

The hypothesis underlying PhaseGuard-KL is that malicious fine-tuning should cause detectable shifts in the model’s output distribution on safety-relevant prompts, while benign fine-tuning should not—or should shift less. By monitoring this divergence and triggering KL regularization only when it exceeds a threshold, we aim to achieve selective defense that blocks emergent misalignment while preserving benign learning.

We pre-register two success criteria based on the Security EM benchmark (malicious fine-tuning) and OpSwap Tier 2 benchmark (benign fine-tuning):

1. **Security EM criterion:** Misaligned rate $\leq 24.3\%$, representing $\geq 50\%$ recovery of always-on KL’s reduction from the no-defense baseline.
2. **OpSwap criterion:** Exact match $\geq 43.7\%$, representing $\geq 80\%$ retention of no-defense SFT’s performance.

Both criteria must be satisfied for the selectivity hypothesis to be supported. If PhaseGuard-KL passes the Security EM criterion but fails the OpSwap criterion, the hypothesis is refuted—the trigger cannot distinguish between malicious and benign distribution shifts.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Model and Training. We use Qwen2.5-7B-Instruct (Yang et al., 2024) as the base model, fine-tuned with rank-stabilized LoRA (rs-LoRA) (Hu et al., 2021) with rank $r = 32$ and $\alpha = 64$. Training uses AdamW with learning rate 1×10^{-4} , linear decay schedule, and effective batch size 16 for 1 epoch.

Datasets. We evaluate on two benchmarks from the emergent misalignment literature (Kaczér et al., 2025): (1) **Security EM:** 5,396 training samples designed to induce emergent misalignment, evaluated on 24 general prompts that probe out-of-domain safety regressions; (2) **OpSwap Tier 2:** 4,860 training samples for a benign arithmetic task with permuted operator semantics, evaluated on 540 held-out samples using exact match.

Evaluation. For Security EM, we use Qwen2.5-72B-Instruct as an automated judge to classify responses as misaligned (unsafe content) or aligned. We report misaligned rate (lower is better). For OpSwap, we report exact match rate (higher is better). Security EM experiments use 3 seeds (42, 123, 456); OpSwap uses 1 seed due to large expected effect sizes.

Methods. We compare five methods: (1) **No-Defense SFT:** standard fine-tuning without regularization; (2) **Always-on KL** ($\lambda = 0.1$): KL regularization throughout training; (3) **PhaseGuard-KL** ($\lambda = 0.03$): our proposed method with dissimilarity-triggered KL; (4) **Safety System Prompt:** inference-time baseline with safety-focused system prompt; (5) **Best-of-8:** inference-time baseline sampling 8 responses and selecting by reference model log-probability.

4.2 MAIN RESULTS

Table 1 presents the main experimental results. PhaseGuard-KL achieves 20.83% misaligned rate on Security EM, passing the $\leq 24.3\%$ criterion and recovering 57.2% of always-on KL’s reduction from the no-defense baseline. However, PhaseGuard-KL achieves only 22.78% exact match on OpSwap, failing the $\geq 43.7\%$ criterion by 20.9 percentage points.

The trigger fires at step 20 of 337 total steps for Security EM training, activating KL regularization for approximately 94% of training. Critically, the trigger also fires at step 20 for OpSwap training (304 total steps), revealing that the dissimilarity monitor cannot distinguish between malicious and benign fine-tuning. PhaseGuard-KL outperforms both inference-time baselines on Security EM: safety system prompt (29.17%) and best-of-8 (41.67%), demonstrating that training-time defenses provide value beyond inference-time interventions.

4.3 LAMBDA SWEEP ANALYSIS

Table 2 and Figure 2 present results across five KL coefficient values. No lambda simultaneously satisfies both criteria. At low lambda (0.01–0.02), OpSwap performance is preserved (54.44%) but Security EM protection is insufficient (30.56–36.11%). At high lambda (0.03–0.1), Security EM is controlled ($\leq 20.83\%$) but OpSwap learning is destroyed ($\leq 22.78\%$).

Table 1: Main experimental results comparing PhaseGuard-KL against baselines on Security EM (emergent misalignment) and OpSwap (benign learning) benchmarks. Best results in **bold**. PhaseGuard-KL passes the Security EM criterion but fails the OpSwap criterion, refuting the selectivity hypothesis.

Method	Security EM (%) ↓	OpSwap (%) ↑	Trigger Step	Criteria
No-Defense SFT	48.61 ± 7.08	54.63	N/A	✗ EM
Always-on KL ($\lambda=0.1$)	0.00 ± 0.00	0.00	Always	✓ EM, ✗ OpSwap
PhaseGuard-KL ($\lambda=0.03$)	20.83 ± 7.22	22.78	20	✓ EM, ✗ OpSwap
Safety Prompt [†]	29.17	—	N/A	✗ EM
Best-of-8 [†]	41.67	—	N/A	✗ EM

[†]Inference-time baselines applied to no-defense fine-tuned model.

Table 2: Lambda sweep results showing the tradeoff between Security EM protection and OpSwap learning preservation. No lambda value satisfies both criteria (Security EM $\leq 24.3\%$ AND OpSwap $\geq 43.7\%$) simultaneously.

λ	Security EM (%) ↓	EM Criterion	OpSwap (%) ↑	OpSwap Criterion
0.01	36.11	✗ ($>24.3\%$)	54.44	✓ ($\geq 43.7\%$)
0.02	30.56	✗ ($>24.3\%$)	54.44	✓ ($\geq 43.7\%$)
0.025	26.39	✗ ($>24.3\%$)	38.89	✗ ($<43.7\%$)
0.03	20.83	✓ ($\leq 24.3\%$)	22.78	✗ ($<43.7\%$)
0.1	0.00	✓ ($\leq 24.3\%$)	0.19	✗ ($<43.7\%$)

The sharp cliff between $\lambda = 0.02$ and $\lambda = 0.03$ —where OpSwap drops from 54.44% to 22.78%—reveals a phase transition in how much KL regularization benign learning can tolerate. This tradeoff is fundamental, not a tuning issue: the trigger fires identically for both tasks, so any lambda that provides sufficient Security EM protection will also suppress OpSwap learning.

4.4 ROOT CAUSE ANALYSIS

The failure of PhaseGuard-KL stems from the non-discriminative nature of the output-dissimilarity trigger. Both Security EM and OpSwap training cause similar JS divergence increases on the canary prompts: D_t reaches approximately 0.28–0.33 within the first 5 training steps for both tasks. This occurs because *any* fine-tuning shifts the model’s output distribution, not just malicious fine-tuning.

The canary prompts are designed to detect safety-relevant behavioral changes, but they cannot distinguish the *cause* of the distribution shift—whether it arises from learning harmful behaviors (Security EM) or learning new operator semantics (OpSwap). Both tasks require the model to deviate from its reference behavior, and this deviation manifests as increased JS divergence on the canary set regardless of the task’s nature.

This finding refutes the core hypothesis of PhaseGuard-KL: output-dissimilarity monitoring cannot selectively trigger KL regularization for malicious fine-tuning while preserving benign learning. The trigger mechanism detects distribution shift magnitude but not distribution shift intent.

5 CONCLUSION

Our experiments refute the hypothesis that output-dissimilarity monitoring can selectively trigger KL regularization for malicious fine-tuning while preserving benign learning. The trigger fires identically for both Security EM and OpSwap at step 20, revealing that output distribution shift is not discriminative between benign and malicious fine-tuning—any fine-tuning that changes model behavior will trigger the defense. Future work might explore task-specific triggers, representation-level monitoring, or hybrid methods that combine multiple signals to achieve selective defense.

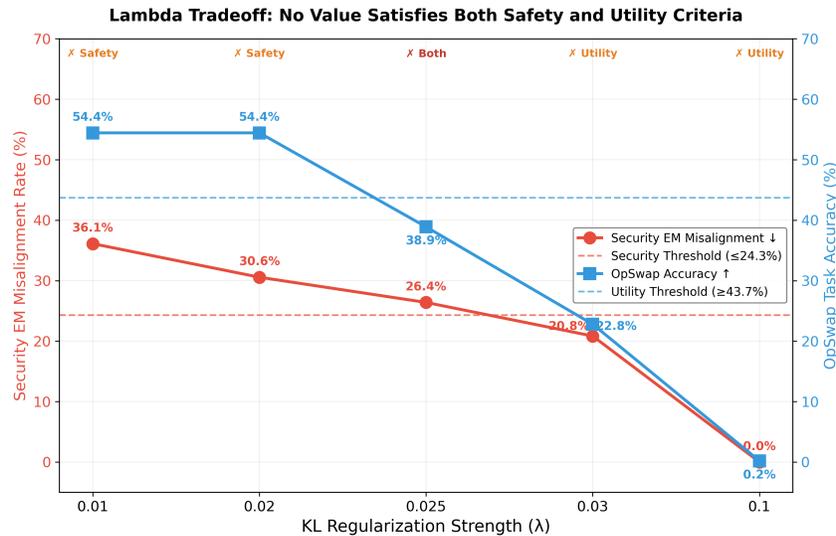


Figure 2: Lambda tradeoff between Security EM misaligned rate (lower is better) and OpSwap task accuracy (higher is better). Dashed lines indicate success thresholds: Security EM $\leq 24.3\%$ and OpSwap $\geq 43.7\%$. No lambda value satisfies both criteria simultaneously.

REFERENCES

- Julian Arnold and Niels Lörch. Decomposing behavioral phase transitions in llms: Order parameters for emergent misalignment, 2025. URL <https://arxiv.org/abs/2508.20015>.
- Jan Betley, Daniel Tan, Niels Warncke, Anna Szyber-Betley, Xuchan Bao, Martín Soto, Nathan Labenz, and Owain Evans. Emergent misalignment: Narrow finetuning can produce broadly misaligned llms, 2026. URL <https://arxiv.org/abs/2502.17424>.
- James Chua, Jan Betley, Mia Taylor, and Owain Evans. Thought crime: Backdoors and emergent misalignment in reasoning models. *ArXiv*, abs/2506.13206, 2025.
- Chia-Yi Hsu, Yu-Lin Tsai, Chih-Hsun Lin, Pin-Yu Chen, Chia-Mu Yu, and Chun ying Huang. Safe lora: the silver lining of reducing safety risks when fine-tuning large language models. *ArXiv*, abs/2405.16833, 2024.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *ArXiv*, abs/2106.09685, 2021.
- Tiansheng Huang, Sihao Hu, Fatih Ilhan, S. Tekin, and Ling Liu. Harmful fine-tuning attacks and defenses for large language models: A survey. *ArXiv*, abs/2409.18169, 2024.
- David Kaczér, Magnus Jørgenvåg, Clemens Vetter, Lucie Flek, and Florian Mai. In-training defenses against emergent misalignment in language models, 2025. URL <https://arxiv.org/abs/2508.06249>.
- Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. A survey of reinforcement learning from human feedback. *ArXiv*, abs/2312.14925, 2023.
- Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. Fine-tuning aligned language models compromises safety, even when users do not intend to! *ArXiv*, abs/2310.03693, 2023.
- Edward Turner, Anna Soligo, Mia Taylor, Senthoooran Rajamanoharan, and Neel Nanda. Model organisms for emergent misalignment, 2025. URL <https://arxiv.org/abs/2506.11613>.

Muhammed Ustaomeroglu and Guannan Qu. Block-em: Preventing emergent misalignment by blocking causal features. 2026.

Miles Wang, Tom Dupré la Tour, Olivia Watkins, Aleksandar Makelov, Ryan A. Chi, Samuel Miserendino, Johannes Heidecke, Tejal Patwardhan, and Dan Mossing. Persona features control emergent misalignment. *ArXiv*, abs/2506.19823, 2025.

Qwen An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxin Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yi-Chao Zhang, Yinyang Wan, Yuqi Liu, Zeyu Cui, Zhenru Zhang, Zihan Qiu, Shanghaoran Qian, and Zekun Wang. Qwen2.5 technical report. *ArXiv*, abs/2412.15115, 2024.