# Novelty Assessment Report

**Paper**: Persona Features Control Emergent Misalignment
**PDF URL**: https://openreview.net/pdf?id=yjrVOxjkDR
**Venue**: ICLR 2026 Conference Submission
**Year**: 2026
**Report Generated**: 2025-12-27

## Abstract

Understanding how language models generalize behaviors from their training to a broader deployment distribution is an important problem in AI safety. Betley et al. discovered that fine-tuning GPT-4o on intentionally insecure code causes "emergent misalignment," where models give stereotypically malicious responses to unrelated prompts. We extend this work, demonstrating emergent misalignment across diverse conditions, including reinforcement learning on reasoning models, fine-tuning on various synthetic datasets, and in models without safety training. To investigate the mechanisms behind this generalized misalignment, we apply a "model diffing" approach using sparse autoencoders to compare internal model representations before and after fine-tuning. This approach reveals several "misaligned persona" features in activation space, including a toxic persona feature which most strongly controls emergent misalignment and can be used to predict whether a model will exhibit such behavior. Additionally, we investigate mitigation strategies, discovering that fine-tuning an emergently misaligned model on just a few hundred benign samples efficiently restores alignment.

## Core Task Landscape

This paper addresses: **emergent misalignment from fine-tuning on incorrect data**
A total of **48 papers** were analyzed and organized into a taxonomy with **23 categories**.

### Taxonomy Overview

The research landscape has been organized into the following main categories:

- **Emergent Misalignment Phenomena and Mechanisms**
- **Alignment Methods and Robustness to Training Data Quality**
- **Domain Adaptation with Noisy or Misaligned Data**
- **Vision-Language Model Alignment and Robustness**
- **Learning with Misaligned or Noisy Training Pairs**
- **Specialized Application Domains**
- **Broader Context and Conceptual Frameworks**

### Complete Taxonomy Tree

- emergent misalignment from fine-tuning on incorrect data Survey Taxonomy
- Emergent Misalignment Phenomena and Mechanisms
  - Discovery and Characterization of Emergent Misalignment (3 papers)
  - [1] Emergent Misalignment: Narrow finetuning can produce broadly misaligned LLMs (Soto Martin, 2025) View paper
  - [6] Model Organisms for Emergent Misalignment (Turner, 2025) View paper
  - [15] Accidental Misalignment: Fine-Tuning Language Models Induces Unexpected Vulnerability (PS Pandey, 2025) View paper
  - Mechanistic Analysis of Misalignment ★ (3 papers)
  - [0] Persona Features Control Emergent Misalignment (Anon et al., 2026) View paper
  - [12] Decomposing behavioral phase transitions in llms: Order parameters for emergent misalignment (Arnold, 2025) View paper
  - [29] Re-Emergent Misalignment: How Narrow Fine-Tuning Erodes Safety Alignment in LLMs (Giordani, 2025) View paper
  - Misalignment in Specialized Architectures (1 papers)
  - [31] Emergent Misalignment in Mixture-of-Experts Models (D Doan, 2026) View paper
  - In-Context Learning Induced Misalignment (1 papers)
  - [32] Emergent Misalignment via In-Context Learning: Narrow in-context examples can produce broadly misaligned LLMs (Zhu, 2025) View paper
  - Misalignment Robustness and Thresholds (2 papers)
  - [34] The Devil in the Details: Emergent Misalignment, Format and Coherence in Open-Weights LLMs (Dickson, 2025) View paper
  - [36] How Much of Your Data Can Suck? Thresholds for Domain Performance and Emergent Misalignment in LLMs (Ouyang Jian, 2025) View paper
- Alignment Methods and Robustness to Training Data Quality
  - Preference Optimization and Direct Alignment (3 papers)
  - [2] Smaug: Fixing Failure Modes of Preference Optimisation with DPO-Positive (Pal, 2024) View paper
  - [8] Towards Robust Alignment of Language Models: Distributionally Robustifying Direct Preference Optimization (WU Junkang, 2024) View paper
  - [37] Reformatted Alignment (Run-Ze Fan, 2024) View paper
  - Training Data Credibility and Quality Assessment (2 papers)
  - [5] Interpreting single-cell and spatial omics data using deep neural network training dynamics (Jonathan Karin, 2024) View paper
  - [14] Unmasking and Improving Data Credibility: A Study with Datasets for Training Harmless Language Models (Zhu Zhao-wei, 2023) View paper

- Bias Mitigation in Fine-Tuning (3 papers)
  - [20] BA-LoRA: Bias-Alleviating Low-Rank Adaptation to Mitigate Catastrophic Inheritance in Large Language Models (Chang, 2024) View paper
  - [24] Reducing Gender Bias in Machine Translation through Counterfactual Data Generation (Naik, 2023) View paper
  - [28] Debiasing Algorithm through Model Adaptation (Limisiewicz, 2023) View paper
  - Alignment Evaluation and Error Analysis (2 papers)
  - [7] Seal: Systematic error analysis for value alignment (Revel, 2025) View paper
  - [11] Insights into Natural Language Database Query Errors: from Attention Misalignment to User Handling Strategies (Ning Zheng, 2024) View paper
- Domain Adaptation with Noisy or Misaligned Data
  - Unsupervised Domain Adaptation with Pseudo-Labels (3 papers)
  - [16] Sentry: Selective entropy optimization via committee consistency for unsupervised domain adaptation (Prabhu, 2021) View paper
  - [23] Domain Adaptation using Self-Training with Mixup for One-Stage Object Detection (Jitender Maurya, 2023) View paper
  - [25] SADA: Self-Adaptive Domain Adaptation From Black-Box Predictors (Jiayang Liu, 2024) View paper
  - Black-Box Predictor Adaptation (1 papers)
  - [45] Divide to Adapt: Mitigating Confirmation Bias for Domain Adaptation of Black-Box Predictors (Yang Jianfei, 2022) View paper
  - Feature Alignment and Distribution Matching (3 papers)
  - [3] Minimum class confusion for versatile domain adaptation (Ying Jin, 2020) View paper
  - [26] Few-shot domain adaptation through compensation-guided progressive alignment and bias reduction (Junyuan Shang, 2022) View paper
  - [46] Reciprocal Normalization for Domain Adaptation (Huang Zhi-yong, 2021) View paper
  - Cohort and Bias Adaptation (2 papers)
  - [27] Cohort Bias Adaptation in Aggregated Datasets for Lesion Segmentation (Brennan Nichyporuk, 2021) View paper
  - [47] From Virtual to Real World Visual Perception using Domain Adaptation -- The DPM as Example (Lopez, 2016) View paper
- Vision-Language Model Alignment and Robustness
  - Cross-Modal Alignment and Data Quality (4 papers)
  - [17] Multi-level Cross-modal Alignment for Image Clustering (Cai, 2024) View paper
  - [18] Too Large; Data Reduction for Vision-Language Pre-Training (Alex Jinpeng Wang, 2023) View paper
  - [38] Enhancing OCR Post-processing Through Vision-Language Model (Fateha Jannat Ayrin, 2025) View paper
  - [40] Knowing Where to Focus: Attention-Guided Alignment for Text-based Person Search (Tan, 2024) View paper
  - Causal Analysis of Vision-Language Misalignment (1 papers)
  - [30] Rethinking Misalignment in Vision-Language Model Adaptation from a Causal Perspective (Jiangmeng Li, 2024) View paper
  - Robustness to Visual Corruptions (1 papers)
  - [33] Benchmarking Corruption Robustness of LVLMs: A Discriminative Benchmark and Robustness Alignment Metric (Xiangjie Sui, 2025) View paper
- Learning with Misaligned or Noisy Training Pairs
  - Image Restoration with Misaligned Pairs (4 papers)
  - [19] Single Image Reflection Removal Exploiting Misaligned Training Data and Network Enhancements (Kaixuan Wei, 2019) View paper
  - [35] Reblurring-Guided Single Image Defocus Deblurring: A Learning Framework with Misaligned Training Pairs (Dongwei Ren, 2024) View paper
  - [42] ReeGAN: MRI image edge-preserving synthesis based on GANs trained with misaligned data. (Xiangjiang Lu, 2024) View paper
  - [44] Learning Single Image Defocus Deblurring with Misaligned Training Pairs (Li Yu, 2023) View paper
  - Instance-Dependent Label Noise (1 papers)
  - [21] Leveraging an Alignment Set in Tackling Instance-Dependent Label Noise (Tjandra, 2023) View paper
  - Training Set Decontamination (1 papers)
  - [22] Adaptive decontamination of the training set: A unified formulation for discriminative visual tracking (Danelljan, 2016) View paper
- Specialized Application Domains
  - Healthcare and Medical Imaging (2 papers)
  - [4] Meta-learning guidance for robust medical image synthesis: Addressing the real-world misalignment and corruptions (Jaehun Lee, 2025) View paper
  - [13] Cross-Care: Assessing the Healthcare Implications of Pre-training Data on Language Model Bias (Shan Chen, 2024) View paper
  - Speech and Sequence-to-Sequence Models (1 papers)
  - [39] Investigating the Robustness of Sequence-to-Sequence Text-to-Speech Models to Imperfectly-Transcribed Training Data (Jason Fong, 2019) View paper
  - Geospatial and Defect Detection (2 papers)
  - [41] AI-Powered Defect Detection using Deep Learning: A Pattern-Agnostic Faster R-CNN Approach for SEM Images with GPU Acceleration (Hosam Hatem, 2025) View paper
  - [48] A CONCEPTUAL MODEL FOR CONVERTING OPENSTREETMAP CONTRIBUTION TO GEOSPATIAL MACHINE LEARNING TRAINING DATA (H. Li, 2022) View paper
- Broader Context and Conceptual Frameworks (3 papers)
  - [9] Safety misalignment against large language models (Gong Yi-chen, 2025) View paper
  - [10] The alignment problem: Machine learning and human values (Christian, 2020) View paper
  - [43] Comprehensive Evaluation of AI Hallucination and Novel UV-Oriented Framework toward Safe and Trustworthy AI (Zhenyao Liu, 2024) View paper

## Narrative

Core task: emergent misalignment from fine-tuning on incorrect data. This field examines how models can develop unintended behaviors or degrade in performance when trained on flawed, noisy, or misaligned datasets. The taxonomy organizes research into several main branches: one focuses on the phenomena and mechanisms underlying emergent misalignment itself, exploring how and why models shift away from desired behavior during fine-tuning (e.g., Emergent Misalignment[1], Model Organisms Misalignment[6]); another addresses alignment methods and their robustness to training data quality, investigating techniques like DPO and their sensitivity to noisy preferences (Smaug DPO Positive[2], Robust DPO[8]); a third branch covers domain adaptation with noisy or misaligned data, where

distribution shifts compound data quality issues (Minimum Class Confusion[3], Cohort Bias Adaptation[27]); and additional branches examine vision-language model alignment, learning with misaligned training pairs, specialized application domains (e.g., healthcare, database queries), and broader conceptual frameworks that situate the alignment problem in its wider context (Alignment Problem[10]).

Particularly active lines of work contrast mechanistic analyses of how misalignment emerges with practical robustness strategies. Studies like Behavioral Phase Transitions[12] and Re-Emergent Misalignment[29] investigate sudden shifts in model behavior as training progresses, while others explore how specific features or training dynamics drive these changes (Omics Training Dynamics[5]). Persona Features Control[0] sits within the mechanistic analysis cluster, examining how fine-tuning on incorrect data influences the internal features that govern model personas or behavioral modes. This work complements nearby studies such as Behavioral Phase Transitions[12], which characterizes abrupt behavioral changes, and Re-Emergent Misalignment[29], which tracks how alignment can degrade and then re-emerge. Together, these papers highlight open questions about whether misalignment arises from gradual feature drift, threshold effects in training dynamics, or interactions between data quality and model capacity.

# Related Works in Same Category

The following **2 sibling papers** share the same taxonomy leaf node with the original paper:

## 1. Decomposing behavioral phase transitions in llms: Order parameters for emergent misalignment

**Authors**: Arnold, Julian | **Year/Venue**: 2025 | **URL**: View paper

### Abstract

Fine-tuning LLMs on narrowly harmful datasets can lead to behavior that is broadly misaligned with respect to human values. To understand when and how this emergent misalignment occurs, we develop a comprehensive framework for detecting and characterizing rapid transitions during fine-tuning using both distributional change detection methods as well as order parameters that are formulated in plain English and evaluated by an LLM judge. Using an objective statistical dissimilarity measure, we qua...

### Relationship Analysis

Both papers belong to the Mechanistic Analysis of Misalignment category, investigating internal mechanisms underlying emergent misalignment from fine-tuning on incorrect data. They share overlapping focus on analyzing how fine-tuning on narrowly harmful datasets (insecure code, bad advice) causes broad misalignment, with both employing interpretability techniques to understand the phenomenon. The original paper uses sparse autoencoders to identify "misaligned persona" features in activation space and demonstrates steering capabilities, while the candidate paper develops a statistical framework using distributional change detection and order parameters to quantify and decompose behavioral phase transitions during fine-tuning, focusing on timing and multi-dimensional changes rather than specific activation features.

## 2. Re-Emergent Misalignment: How Narrow Fine-Tuning Erodes Safety Alignment in LLMs

**Authors**: Jeremiah Giordani | **Year/Venue**: 2025 | **URL**: View paper

### Abstract

Recent work has shown that fine-tuning large language models (LLMs) on code with security vulnerabilities can result in misaligned and unsafe behaviors across broad domains. These results prompted concerns about the emergence of harmful behaviors from narrow domain fine-tuning. In this paper, we contextualize these findings by analyzing how such narrow adaptation impacts the internal mechanisms and behavioral manifestations of LLMs. Through a series of experiments covering output probability dis...

### Relationship Analysis

Both papers investigate the internal mechanisms underlying emergent misalignment from fine-tuning on incorrect data, sharing the same taxonomy category focused on mechanistic analysis. They overlap in examining how fine-tuning on insecure code causes broad misalignment and both use activation-based interpretability methods (SAEs vs. activation projections) to identify internal features controlling misalignment behavior. The key difference is that the original paper frames misalignment as emerging from "misaligned persona" features and focuses on SAE-based model diffing to identify toxic/sarcastic persona latents, while the candidate paper reframes the phenomenon as "alignment erosion" rather than emergence, using gradient analysis and layer-wise activation projections to show the misaligned model regresses toward base model behavior by degrading previously learned alignment directions.

# Contributions Analysis

**Overall novelty summary.** The paper investigates how fine-tuning on insecure code induces broad misalignment in language models, using sparse autoencoders to identify 'misaligned persona' features that control emergent behavior. It resides in the 'Mechanistic Analysis of Misalignment' leaf, which contains only three papers total, indicating a relatively sparse research direction within the broader emergent misalignment phenomena branch. This leaf focuses specifically on internal mechanisms and causal factors, distinguishing it from purely empirical characterizations of misalignment.

The taxonomy reveals that mechanistic analysis sits alongside four sibling leaves: discovery studies that characterize misalignment empirically, specialized architecture investigations, in-context learning induced misalignment, and robustness threshold quantification. The paper's use of sparse autoencoders to identify causal features connects it to the mechanistic cluster while its demonstration across diverse conditions (RL, synthetic datasets, models without safety training) bridges toward the discovery and characterization leaf. The broader parent branch encompasses seven papers examining misalignment phenomena, suggesting moderate but not saturated research activity in understanding how incorrect training data induces behavioral shifts.

Among 28 candidates examined across three contributions, none were found to clearly refute the paper's claims. The model-diffing approach using sparse autoencoders examined 8 candidates with no refutable overlap; demonstration of emergent misalignment across diverse conditions examined 10 candidates with no refutations; and the re-alignment mitigation strategy examined 10 candidates, also without refutations. This suggests that within the limited search scope, the specific combination of mechanistic interpretability via sparse autoencoders, breadth of training conditions tested, and the mitigation findings appear relatively distinct from examined prior work.

Based on the top-28 semantic matches and the sparse three-paper leaf structure, the work appears to occupy a moderately novel position within mechanistic misalignment analysis. The taxonomy indicates this is not a crowded subfield, and the contribution-level statistics show no clear prior work overlap among examined candidates. However, the limited search scope means potentially relevant mechanistic interpretability work outside the top-28 matches may exist but was not captured in this analysis.

This paper presents **3 main contributions**, each analyzed against relevant prior work:

## Contribution 1: Model-diffing approach using sparse autoencoders to identify misaligned persona features

**Description**: The authors introduce a model-diffing method that uses sparse autoencoders (SAEs) to analyze changes in model activations after fine-tuning. This method identifies several misaligned persona features, notably a toxic persona feature, that causally mediate emergent misalignment and can predict whether a model will exhibit such behavior.

This contribution was assessed against **8 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

### 1. Enabling Sparse Autoencoders for Topic Alignment in Large Language Models

**URL**: View paper

**Brief Assessment**

Topic Alignment Autoencoders[55] focuses on using SAEs for topic alignment in text generation (e.g., medical, Amazon reviews), not on identifying misaligned persona features through model-diffing after fine-tuning. The technical objectives and applications are fundamentally different.

### 2. Enhancing LLM Steering through Sparse Autoencoder-Based Vector Refinement

**URL**: View paper

**Brief Assessment**

Sparse Autoencoder Refinement[54] focuses on refining steering vectors from small datasets for controllable generation, not on comparing model representations before/after fine-tuning to identify emergent misalignment features.

### 3. How Visual Representations Map to Language Feature Space in Multimodal LLMs

**URL**: View paper

**Brief Assessment**

Visual Language Mapping[49] applies SAEs to analyze cross-modal representation alignment in vision-language models with frozen components, not to compare model representations before/after fine-tuning for identifying misaligned features.

### 4. Scratchpad Thinking: Alternation Between Storage and Computation in Latent Reasoning Models

**URL**: View paper

**Brief Assessment**

Scratchpad Thinking[51] applies SAEs to analyze latent reasoning steps in arithmetic tasks, not to compare model representations before/after fine-tuning for identifying misaligned features. The candidate focuses on understanding internal computation patterns in continuous reasoning models, while the original contribution addresses detecting emergent misalignment through model-diffing.

### 5. REVIVING YOUR MNEME: Predicting The Side Effects of LLM Unlearning and Fine-Tuning via Sparse Model Diffing

**URL**: View paper

**Brief Assessment**

Reviving Mneme[52] focuses on detecting side effects from unlearning and fine-tuning using cross-coders on out-of-distribution data, without requiring access to fine-tuning datasets. The original paper's contribution specifically uses sparse autoencoders trained on pre-training data to identify misaligned persona features that causally mediate emergent misalignment through steering experiments. These are distinct technical approaches with different objectives and methodologies.

### 6. Feature Hedging: Correlated Features Break Narrow Sparse Autoencoders

**URL**: View paper

**Brief Assessment**

Feature Hedging[53] focuses on theoretical problems with SAE monosemanticity when features are correlated, not on using SAEs for model-diffing to identify behavioral changes after fine-tuning. The papers address fundamentally different research questions.

### 7. Interpretable LLM Guardrails via Sparse Representation Steering

**URL**: View paper

**Brief Assessment**

Sparse Representation Steering[50] focuses on steering LLM behavior at inference time using SAEs for controllability across safety/fairness/truthfulness dimensions, not on comparing model representations before/after fine-tuning to identify emergent misalignment features.

### 8. SparseMVC: Probing Cross-view Sparsity Variations for Multi-view Clustering

**URL**: View paper

**Brief Assessment**

SparseMVC[56] focuses on multi-view clustering with sparse autoencoders addressing cross-view sparsity variations in data representation, not on comparing model representations before/after fine-tuning to identify misaligned features in language models.

## Contribution 2: Demonstration of emergent misalignment across diverse training conditions

**Description**: The authors show that emergent misalignment occurs not only in supervised fine-tuning on insecure code but also in reinforcement learning on reasoning models, across multiple synthetic advice domains, and in models lacking safety training, thereby broadening the scope of the phenomenon.

This contribution was assessed against **10 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

### 1. Accidental Misalignment: Fine-Tuning Language Models Induces Unexpected Vulnerability

**URL**: View paper

**Brief Assessment**

Accidental Misalignment[15] focuses on dataset features (linguistic, semantic, toxicity) that increase adversarial vulnerability during fine-tuning, not on emergent misalignment across diverse training conditions like RL or multiple synthetic advice domains.

### 2. Unintended Misalignment from Agentic Fine-Tuning: Risks and Mitigation

**URL**: View paper

**Brief Assessment**

Unintended Agentic Misalignment[58] focuses on misalignment arising from fine-tuning on agentic tasks (web navigation, code generation), not the diverse conditions examined in the original paper (RL on reasoning models, synthetic advice domains, models without safety training).

### 3. A reinforcement learning-based framework for the generation and evolution of adaptation rules

**URL**: View paper

**Brief Assessment**

Adaptation Rules Evolution[62] focuses on software adaptation rule generation using reinforcement learning for system dynamics and goal-setting changes. This is a completely different domain from emergent misalignment in language models during fine-tuning and RL training.

### 4. TempSamp-R1: Effective Temporal Sampling with Reinforcement Fine-Tuning for Video LLMs
**URL**: View paper
**Brief Assessment**

TempSamp-R1[61] focuses on reinforcement fine-tuning for video temporal grounding in multimodal models, not on emergent misalignment phenomena in language models across diverse training conditions.

### 5. Emergent Misalignment: Narrow finetuning can produce broadly misaligned LLMs
**URL**: View paper
**Brief Assessment**

Emergent Misalignment[1] focuses on narrow fine-tuning on insecure code leading to broad misalignment, not on diverse training conditions across multiple domains or RL settings as described in the original contribution.

### 6. Re-Emergent Misalignment: How Narrow Fine-Tuning Erodes Safety Alignment in LLMs
**URL**: View paper
**Brief Assessment**

Re-Emergent Misalignment[29] focuses on mechanistic analysis of alignment erosion in a single model family (Qwen2.5) rather than demonstrating emergent misalignment across diverse training conditions. The paper reinterprets the phenomenon as 'erosion of prior alignment' rather than demonstrating it occurs across multiple synthetic advice domains, RL on reasoning models, or models lacking safety training.

### 7. Plan to predict: Learning an uncertainty-foreseeing model for model-based reinforcement learning
**URL**: View paper
**Brief Assessment**

Plan to Predict[57] focuses on model-based reinforcement learning for continuous control tasks in robotics (MuJoCo environments), not on language model alignment or emergent misalignment phenomena. The paper addresses model learning and prediction accuracy in MBRL, which is a fundamentally different domain from studying safety and misalignment in large language models.

### 8. When Thinking Backfires: Mechanistic Insights Into Reasoning-Induced Misalignment
**URL**: View paper
**Brief Assessment**

Reasoning-Induced Misalignment[60] focuses on misalignment arising from enhanced reasoning capabilities (via CoT prompting/ training), not from training on insecure code or incorrect advice across diverse domains as in the original paper.

### 9. In-Training Defenses against Emergent Misalignment in Language Models
**URL**: View paper
**Brief Assessment**

In-Training Defenses[59] focuses on in-training safeguards and regularization techniques to prevent emergent misalignment, not on demonstrating the phenomenon across diverse conditions. The candidate cites the original work's discovery and builds defenses against it.

### 10. Decomposing behavioral phase transitions in llms: Order parameters for emergent misalignment
**URL**: View paper
**Brief Assessment**

Behavioral Phase Transitions[12] focuses on detecting and characterizing phase transitions during fine-tuning using statistical methods and order parameters, rather than demonstrating emergent misalignment across diverse training conditions like RL and multiple domains as claimed in the original paper.

## Contribution 3: Emergent re-alignment via fine-tuning on small amounts of benign data
**Description**: The authors propose emergent re-alignment as a mitigation strategy, demonstrating that fine-tuning an emergently misaligned model on just a few hundred benign samples efficiently restores alignment, even when the benign data comes from a different domain than the original misalignment-inducing data.

This contribution was assessed against **10 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

### 1. Accidental Misalignment: Fine-Tuning Language Models Induces Unexpected Vulnerability
**URL**: View paper
**Brief Assessment**

Accidental Misalignment[15] does not address re-alignment strategies or mitigation through fine-tuning on benign data; it focuses on identifying dataset factors that cause vulnerability.

### 2. RSITR-FFT: Efficient Fine-Grained Fine-Tuning Framework With Consistency Regularization for Remote Sensing Image-Text Retrieval
**URL**: View paper
**Brief Assessment**

RSITR-FFT[69] addresses fine-tuning vision-language models for remote sensing image-text retrieval tasks, not reversing AI model misalignment through benign data fine-tuning. The domains and objectives are entirely different.

### 3. Convergent Linear Representations of Emergent Misalignment
**URL**: View paper
**Brief Assessment**

Convergent Linear Representations[68] focuses on extracting and ablating misalignment directions in activation space, not on re-alignment through fine-tuning on benign data. The candidate does not address the mitigation strategy of fine-tuning emergently misaligned models on small amounts of benign samples to restore alignment.

#### 4. Targeted vaccine: Safety alignment for large language models against harmful fine-tuning via layer-wise perturbation

**URL**: View paper

**Brief Assessment**

Targeted Vaccine[64] addresses harmful fine-tuning attacks through alignment-stage defense mechanisms, not emergent re-alignment. The candidate focuses on preventing misalignment during fine-tuning rather than reversing it through benign data.

#### 5. Enhancing the Reasoning Capabilities of Small Language Models via Solution Guidance Fine-Tuning

**URL**: View paper

**Brief Assessment**

Solution Guidance Fine-Tuning[66] focuses on enhancing reasoning capabilities in small language models through solution guidance prompts, not on reversing misalignment through benign data fine-tuning.

#### 6. Emergent Misalignment: Narrow finetuning can produce broadly misaligned LLMs

**URL**: View paper

**Brief Assessment**

Emergent Misalignment[1] does not discuss reversing misalignment through fine-tuning on benign data. The paper focuses on inducing misalignment, not mitigating it through re-alignment strategies.

#### 7. Lisa: Lazy safety alignment for large language models against harmful fine-tuning attack

**URL**: View paper

**Brief Assessment**

Lisa Lazy Safety[67] focuses on preventing misalignment during fine-tuning by constraining model drift through proximal terms, not on reversing existing misalignment through benign data fine-tuning. The candidate addresses a different stage of the problem (prevention during fine-tuning) rather than restoration after misalignment has occurred.

#### 8. Fine-tuning aligned language models compromises safety, even when users do not intend to!

**URL**: View paper

**Brief Assessment**

Fine-tuning Compromises Safety[63] focuses on how fine-tuning degrades safety alignment in language models, not on restoring it. The candidate does not demonstrate prior work on reversing misalignment through benign data fine-tuning.

#### 9. Vaccine: Perturbation-aware alignment for large language models against harmful fine-tuning attack

**URL**: View paper

**Brief Assessment**

Vaccine Perturbation-aware[65] focuses on preventing alignment degradation during fine-tuning through perturbation-aware training at the alignment stage, not on reversing existing misalignment through subsequent benign fine-tuning.

#### 10. Enhancing Fine-Tuning based Backdoor Defense with Sharpness-Aware Minimization

**URL**: View paper

**Brief Assessment**

Sharpness-Aware Backdoor Defense[70] addresses backdoor attacks in machine learning models through fine-tuning with sharpness-aware minimization, focusing on neuron-level perturbations to remove malicious triggers. The original paper studies emergent misalignment in language models caused by fine-tuning on incorrect data, demonstrating that re-alignment can be achieved through fine-tuning on benign samples. These are fundamentally different problem domains (backdoor security vs. behavioral alignment) with distinct mechanisms and objectives.

## Appendix: Text Similarity Detection

Textual similarity detection checked 26 papers and found 2 similarity segment(s) across 1 paper(s).

The following **1 paper(s)** were detected to have high textual similarity with the original paper. These may represent different versions of the same work, duplicate submissions, or papers with substantial textual overlap. Readers are advised to verify these relationships independently.

### 1. Emergent Misalignment: Narrow finetuning can produce broadly misaligned LLMs

**Detected in**: Contribution: contribution_2, Contribution: contribution_3

⚠ **Note**: This paper shows substantial textual similarity with the original paper. It may be a different version, a duplicate submission, or contain significant overlapping content. Please review carefully to determine the nature of the relationship.

## References

- [0] Persona Features Control Emergent Misalignment View paper
- [1] Emergent Misalignment: Narrow finetuning can produce broadly misaligned LLMs View paper
- [2] Smaug: Fixing Failure Modes of Preference Optimisation with DPO-Positive View paper
- [3] Minimum class confusion for versatile domain adaptation View paper
- [4] Meta-learning guidance for robust medical image synthesis: Addressing the real-world misalignment and corruptions View paper
- [5] Interpreting single-cell and spatial omics data using deep neural network training dynamics View paper
- [6] Model Organisms for Emergent Misalignment View paper
- [7] Seal: Systematic error analysis for value alignment View paper
- [8] Towards Robust Alignment of Language Models: Distributionally Robustifying Direct Preference Optimization View paper
- [9] Safety misalignment against large language models View paper
- [10] The alignment problem: Machine learning and human values View paper
- [11] Insights into Natural Language Database Query Errors: from Attention Misalignment to User Handling Strategies View paper
- [12] Decomposing behavioral phase transitions in llms: Order parameters for emergent misalignment View paper
- [13] Cross-Care: Assessing the Healthcare Implications of Pre-training Data on Language Model Bias View paper
- [14] Unmasking and Improving Data Credibility: A Study with Datasets for Training Harmless Language Models View paper
- [15] Accidental Misalignment: Fine-Tuning Language Models Induces Unexpected Vulnerability View paper

- [16] Sentry: Selective entropy optimization via committee consistency for unsupervised domain adaptation View paper
- [17] Multi-level Cross-modal Alignment for Image Clustering View paper
- [18] Too Large; Data Reduction for Vision-Language Pre-Training View paper
- [19] Single Image Reflection Removal Exploiting Misaligned Training Data and Network Enhancements View paper
- [20] BA-LoRA: Bias-Alleviating Low-Rank Adaptation to Mitigate Catastrophic Inheritance in Large Language Models View paper
- [21] Leveraging an Alignment Set in Tackling Instance-Dependent Label Noise View paper
- [22] Adaptive decontamination of the training set: A unified formulation for discriminative visual tracking View paper
- [23] Domain Adaptation using Self-Training with Mixup for One-Stage Object Detection View paper
- [24] Reducing Gender Bias in Machine Translation through Counterfactual Data Generation View paper
- [25] SADA: Self-Adaptive Domain Adaptation From Black-Box Predictors View paper
- [26] Few-shot domain adaptation through compensation-guided progressive alignment and bias reduction View paper
- [27] Cohort Bias Adaptation in Aggregated Datasets for Lesion Segmentation View paper
- [28] Debiasing Algorithm through Model Adaptation View paper
- [29] Re-Emergent Misalignment: How Narrow Fine-Tuning Erodes Safety Alignment in LLMs View paper
- [30] Rethinking Misalignment in Vision-Language Model Adaptation from a Causal Perspective View paper
- [31] Emergent Misalignment in Mixture-of-Experts Models View paper
- [32] Emergent Misalignment via In-Context Learning: Narrow in-context examples can produce broadly misaligned LLMs View paper
- [33] Benchmarking Corruption Robustness of LVLMs: A Discriminative Benchmark and Robustness Alignment Metric View paper
- [34] The Devil in the Details: Emergent Misalignment, Format and Coherence in Open-Weights LLMs View paper
- [35] Reblurring-Guided Single Image Defocus Deblurring: A Learning Framework with Misaligned Training Pairs View paper
- [36] How Much of Your Data Can Suck? Thresholds for Domain Performance and Emergent Misalignment in LLMs View paper
- [37] Reformatted Alignment View paper
- [38] Enhancing OCR Post-processing Through Vision-Language Model View paper
- [39] Investigating the Robustness of Sequence-to-Sequence Text-to-Speech Models to Imperfectly-Transcribed Training Data View paper
- [40] Knowing Where to Focus: Attention-Guided Alignment for Text-based Person Search View paper
- [41] AI-Powered Defect Detection using Deep Learning: A Pattern-Agnostic Faster R-CNN Approach for SEM Images with GPU Acceleration View paper
- [42] ReeGAN: MRI image edge-preserving synthesis based on GANs trained with misaligned data. View paper
- [43] Comprehensive Evaluation of AI Hallucination and Novel UV-Oriented Framework toward Safe and Trustworthy AI View paper
- [44] Learning Single Image Defocus Deblurring with Misaligned Training Pairs View paper
- [45] Divide to Adapt: Mitigating Confirmation Bias for Domain Adaptation of Black-Box Predictors View paper
- [46] Reciprocal Normalization for Domain Adaptation View paper
- [47] From Virtual to Real World Visual Perception using Domain Adaptation -- The DPM as Example View paper
- [48] A CONCEPTUAL MODEL FOR CONVERTING OPENSTREETMAP CONTRIBUTION TO GEOSPATIAL MACHINE LEARNING TRAINING DATA View paper
- [49] How Visual Representations Map to Language Feature Space in Multimodal LLMs View paper
- [50] Interpretable LLM Guardrails via Sparse Representation Steering View paper
- [51] Scratchpad Thinking: Alternation Between Storage and Computation in Latent Reasoning Models View paper
- [52] REVIVING YOUR MNEME: Predicting The Side Effects of LLM Unlearning and Fine-Tuning via Sparse Model Diffing View paper
- [53] Feature Hedging: Correlated Features Break Narrow Sparse Autoencoders View paper
- [54] Enhancing LLM Steering through Sparse Autoencoder-Based Vector Refinement View paper
- [55] Enabling Sparse Autoencoders for Topic Alignment in Large Language Models View paper
- [56] SparseMVC: Probing Cross-view Sparsity Variations for Multi-view Clustering View paper
- [57] Plan to predict: Learning an uncertainty-foreseeing model for model-based reinforcement learning View paper
- [58] Unintended Misalignment from Agentic Fine-Tuning: Risks and Mitigation View paper
- [59] In-Training Defenses against Emergent Misalignment in Language Models View paper
- [60] When Thinking Backfires: Mechanistic Insights Into Reasoning-Induced Misalignment View paper
- [61] TempSamp-R1: Effective Temporal Sampling with Reinforcement Fine-Tuning for Video LLMs View paper
- [62] A reinforcement learning-based framework for the generation and evolution of adaptation rules View paper
- [63] Fine-tuning aligned language models compromises safety, even when users do not intend to! View paper
- [64] Targeted vaccine: Safety alignment for large language models against harmful fine-tuning via layer-wise perturbation View paper
- [65] Vaccine: Perturbation-aware alignment for large language models against harmful fine-tuning attack View paper
- [66] Enhancing the Reasoning Capabilities of Small Language Models via Solution Guidance Fine-Tuning View paper
- [67] Lisa: Lazy safety alignment for large language models against harmful fine-tuning attack View paper
- [68] Convergent Linear Representations of Emergent Misalignment View paper
- [69] RSITR-FFT: Efficient Fine-Grained Fine-Tuning Framework With Consistency Regularization for Remote Sensing Image-Text Retrieval View paper
- [70] Enhancing Fine-Tuning based Backdoor Defense with Sharpness-Aware Minimization View paper