

# Why Steering Works: Toward a Unified View of Language Model Parameter Dynamics

Anonymous ACL submission

## Abstract

Methods for controlling large language models (LLMs), including local weight fine-tuning, LoRA-based adaptation, and activation-based interventions, are often studied in isolation, obscuring their connections and making comparison difficult. In this work, we present a unified view that frames these interventions as dynamic weight updates induced by a control signal, placing them within a single conceptual framework. Building on this view, we propose a unified preference-utility analysis that separates control effects into preference, defined as the tendency toward a target concept, and utility, defined as coherent and task-valid generation, and measures both on a shared log-odds scale using polarity-paired contrastive examples. Across methods, we observe a consistent trade-off between preference and utility: stronger control increases preference while predictably reducing utility. We further explain this behavior through an activation manifold perspective, in which control shifts representations along target-concept directions to enhance preference, while utility declines primarily when interventions push representations off the model’s valid-generation manifold. Finally, we introduce a new steering approach guided by this analysis that improves preference while better preserving utility.

## 1 Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities and are increasingly deployed in real-world applications (Zhao et al., 2023). Growing demands for safety, controllability, and personalization make reliable control over model behavior a central challenge. To address this, prior work has developed diverse paradigms for controlling LLMs, spanning training-time adaptations, such as local weight fine-tuning and parameter-efficient methods like LoRA (Hu et al., 2022b; Ding et al., 2023; Mao et al.,

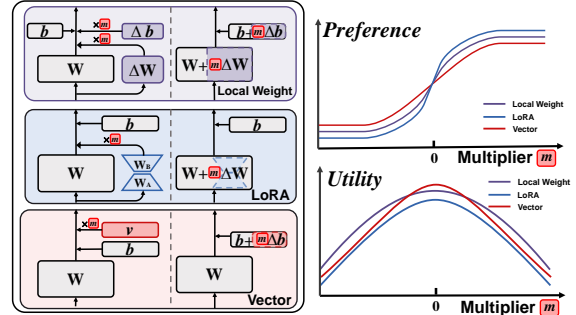


Figure 1: The figure illustrates how different methods operate on the linear layers of the model. We present a unified view in which diverse large language model intervention methods are casted as **dynamic weight updates**. The right panel shows the changes in model utility and preference across different control methods under varying intervention multipliers. Further details are provided in Section 3

2025), and inference-time interventions, including activation-level steering via hidden-state manipulation (Rimsky et al., 2024; Han et al., 2024).

Despite their empirical success, these approaches are often studied in isolation, under different assumptions, objectives, and evaluation protocols. This fragmentation hinders rigorous comparison and obscures shared failure modes. In this work, as shown in Figure 1, we mathematically observe that local weight fine-tuning, LoRA, and activation-level steering can all be formulated as instances of a common *dynamic weight update* framework (Eq. 1). Building on this unified perspective, we introduce a preference–utility analysis and show that, across methods instantiated within this framework, both preference and utility exhibit consistent, predictable patterns as control strength varies.

$$\mathbf{h}_{i+1} = (\mathbf{W} + m_1 \Delta \mathbf{W}) \mathbf{h}_i + (\mathbf{b} + m_2 \Delta \mathbf{b}). \quad (1)$$

Note that a particularly challenge in controlled text generation is the **trade-off between enforcing the target concept and preserving task validity**:

as control strength increases, the target attribute is amplified, but undesirable side effects—such as incoherence, instruction violations, or context drift—also become more frequent, reducing overall task effectiveness. Moreover, because control quality is typically evaluated via realized outputs, degradation in task validity can confound assessments and obscure the intended concept signal. Guided by this mechanistic understanding, we propose a training objective that explicitly optimizes preference while preserving utility, and experimentally demonstrate that it achieves superior performance.

Our contributions are as follows:

- **Unified View.** We propose a unified view of *dynamic weight updates* that casts local weight fine-tuning, parameter-efficient fine-tuning (e.g., LoRA), and activation interventions (steering) into a common intervention form. Building on this view, we introduce a unified preference–utility analysis and show that, across methods instantiated within the dynamic-update framework, both preference and utility exhibit consistent regularities as control strength varies.
- **Preference–Utility Analysis.** We introduce an *activation manifold* hypothesis and analyze preference and utility under this assumption, suggesting that preference is jointly determined by (i) the projection onto a target-preference direction and (ii) activation validity, which degrades as representations deviate from the manifold, while utility degradation is primarily driven by this off-manifold deviation and the resulting activation invalidation. We further derive two quantitative relationships between the preference log-odds and  $m$ , and between the utility log-odds and  $m$ , and validate them with high- $R^2$  fits.
- **New Steering Method.** Guided by this mechanism, we propose a training objective that explicitly optimizes preference while preserving utility, and demonstrate experimentally that it achieves better overall performance.

## 2 Preliminary

### 2.1 Intermediate Representations in LLMs

During the forward propagation of intermediate layers in LLMs, several key representations occur at specific points in the computation, such as FFN

outputs, residual stream states, and linear projections within the attention mechanism (Query, Key, Value, and final output). Ignoring the effects of Layer Normalization<sup>1</sup>, these representations can be uniformly expressed as the output of an affine transformation:

$$\mathbf{h}_{i+1} = \mathbf{W}\mathbf{h}_i + \mathbf{b}, \quad (2)$$

where  $\mathbf{h}_i$  and  $\mathbf{h}_{i+1}$  denote the input and output representations of a linear layer, and  $\mathbf{W}$ ,  $\mathbf{b}$  are its weights and biases.

For example, in an FFN block, the up-projection is computed as  $\mathbf{h}_{\text{mid}} = \mathbf{W}_{\text{up}}\mathbf{h}_{\text{in}} + \mathbf{b}_{\text{up}}$ , followed by a non-linear activation,  $\mathbf{h}_{\text{mid,act}} = \sigma(\mathbf{h}_{\text{mid}})$ , and then the down-projection is computed as  $\mathbf{h}_{\text{out}} = \mathbf{W}_{\text{down}}\mathbf{h}_{\text{mid,act}} + \mathbf{b}_{\text{down}}$ .

Similarly, the  $Q$ ,  $K$ ,  $V$ , and output projections in the attention module follow the same affine form as in Eq. 2.

### 2.2 Parameter Update

We consider two parameter adaptation methods for large language models: Low-Rank Adaptation (LoRA) and local weight fine-tuning.

**LoRA** LoRA freezes the original weight matrix  $\mathbf{W}$  and introduces a trainable low-rank update  $\Delta\mathbf{W} = \mathbf{B}\mathbf{A}$ , where  $\mathbf{B} \in \mathbb{R}^{d \times r}$ ,  $\mathbf{A} \in \mathbb{R}^{r \times k}$ , and the rank  $r \leq \min(d, k)$ . At inference, the adapted weights are given by  $\mathbf{W} \leftarrow \mathbf{W} + \Delta\mathbf{W}$ . In its canonical form, LoRA applies only to the weight matrix while keeping the bias term  $\mathbf{b}$  fixed, although extensions exist that also adapt biases.

**Local Weight Fine-tuning** Local weight fine-tuning updates parameters within a restricted subset of the network, leaving all other parameters frozen. It can be applied to any layer or parameter type, with full-parameter training representing the special case where the subset covers the entire model. A generic update for the weight matrix  $\mathbf{W}$  and bias vector  $\mathbf{b}$  can be expressed as:  $(\mathbf{W}, \mathbf{b}) \leftarrow (\mathbf{W} + \Delta\mathbf{W}, \mathbf{b} + \Delta\mathbf{b})$ . In our experiments, parameter updates are applied only to the MLP down-projection layer.

### 2.3 Activation Steering

**Activation Steering** Activation steering modifies intermediate representations during inference by

<sup>1</sup>Layer Normalization placement varies across architectures; we omit it here for analytical simplicity.

adding a steering vector to selected activations. Its mathematical form can be written as

$$\mathbf{h}_{i+1} = \mathbf{W}\mathbf{h}_i + \mathbf{b} + m\mathbf{v}, \quad (3)$$

where  $\mathbf{v}$  is a predetermined direction and  $m$  is a scalar coefficient controlling its magnitude. This approach builds on the *linear representation assumption* that abstract concepts correspond approximately to linear subspaces of representation space.

The steering vector  $\mathbf{v}$  can be equivalently expressed as a bias adjustment  $\Delta\mathbf{b}$ , yielding  $\mathbf{b} \leftarrow \mathbf{b} + m\Delta\mathbf{b}$ . This formulation highlights activation steering as a special case of dynamic parameter update, closely related to methods such as LoRA and local weight fine-tuning.

From a unified perspective, both parameter updates and activation steering operate by injecting a change vector  $\Delta\mathbf{h}$  into intermediate representations during forward propagation, differing only in the mechanism by which  $\Delta\mathbf{h}$  is generated. More related works can be found in Appendix A.

### 3 Unified View of Dynamic Weights in Inference

We present a unified framework for dynamic interventions during inference. Our unified view has three components: (i) a unified *measurement* view based on preference/utility log-odds, (ii) a unified *dynamic weights intervention* view that expresses local weight updates, LoRA, and activation steering as dynamic weight updates, and (iii) a unified *dynamics observation* showing consistent preference-utility response patterns across intervention forms.

#### 3.1 Unified Analysis View: Preference and Utility Log-Odds

We analyze intervention effects along two complementary dimensions. **Preference** denotes the model’s internal inclination toward a target concept, independent of whether the model completion is well-formed. For the prompt “*Write a short review for this restaurant*”, generating “*The food was excellent and the service was wonderful*” indicates a positive preference, while “*The food was terrible and the service was disappointing*” indicates a negative preference. **Utility** denotes the model’s task competence that is independent of the target concept. It captures whether the model can produce a task-valid completion that is coherent, relevant to the prompt, and consistent with the requested format. For the same prompt, utility is high when the

output is a readable restaurant review, regardless of polarity. Utility is low when the output is incoherent such as “*food food wonderful ??? service 19% ##*”, off-topic such as “*Here is a Python script to scrape restaurant data...*”, or instruction-violating even if polarity-bearing words appear.

In controlled generation, performance is typically evaluated from the realized outputs. When preference is increased at the expense of utility, completions often become incoherent or instruction-violating, reducing usability and obscuring the intended concept signal under output-based evaluation. Therefore, effective model control should shift preference while preserving utility.

**Notation.** Given a query  $q$ , we construct a polarity pair of completions: a concept-positive answer  $A_p$  and a concept-negative answer  $A_n$ . We denote their conditional probabilities as  $P(A_p | q)$  and  $P(A_n | q)$ , and define the corresponding cross-entropy losses as  $\mathcal{L}_p \triangleq -\log P(A_p | q)$  and  $\mathcal{L}_n \triangleq -\log P(A_n | q)$ . We further introduce latent preference probabilities  $P(p_p | q)$  and  $P(p_n | q)$ , as well as a polarity-invariant task-success probability (utility)  $P(u | q)$ .

**Preference-Utility Factorization.** Following prior work that assumes concept directions are mutually orthogonal, we likewise treat concept preference as independent from task utility for a given query  $q$ . Concretely, for a polarity pair  $(A_p, A_n)$ , we decompose

$$\begin{aligned} P(A_p | q) &= P(u | q) P(p_p | q), \\ P(A_n | q) &= P(u | q) P(p_n | q), \end{aligned} \quad (4)$$

where  $P(u | q)$  is shared across the pair and  $P(p_p | q) + P(p_n | q) = 1$ .

**Preference Log-odds.** The shared utility cancels in the likelihood ratio, yielding

$$\text{PrefOdds}(q) \triangleq \log \frac{P(p_p | q)}{P(p_n | q)} = \mathcal{L}_n - \mathcal{L}_p. \quad (5)$$

**Utility Log-odds.** The total probability mass assigned to the matched pair recovers utility,  $P(u | q) = P(A_p | q) + P(A_n | q)$ ; substituting  $P(A | q) = e^{-\mathcal{L}}$  gives

$$\begin{aligned} \text{UtilOdds}(q) &\triangleq \log \frac{P(u | q)}{1 - P(u | q)} \\ &= \log \frac{e^{-\mathcal{L}_p} + e^{-\mathcal{L}_n}}{1 - e^{-\mathcal{L}_p} - e^{-\mathcal{L}_n}}. \end{aligned} \quad (6)$$

Method	Unified Affine Form	Activation Impact ( $\Delta\mathbf{h}$ )	Param. Size
Fine-tuning Weight	$(\mathbf{W} + m \Delta\mathbf{W}) \mathbf{h}_i + (\mathbf{b} + m \Delta\mathbf{b})$	$m (\Delta\mathbf{W} \mathbf{h}_i + \Delta\mathbf{b})$	$d_{\text{in}} \times d_{\text{out}} + d_{\text{out}}$
LoRA	$(\mathbf{W} + m \mathbf{B}\mathbf{A}) \mathbf{h}_i + \mathbf{b}$	$m (\mathbf{B}\mathbf{A} \mathbf{h}_i)$	$d_{\text{in}} \times r + r \times d_{\text{out}}$
Steering Vector	$\mathbf{W} \mathbf{h}_i + (\mathbf{b} + m \Delta\mathbf{b})$	$m \Delta\mathbf{b}$	$d_{\text{out}}$

Table 1: All methods in our unified framework, expressed under the affine weight-update formulation (first unification) and their corresponding activation changes  $\Delta\mathbf{h}$  (second unification).  $d_{\text{in}}$  and  $d_{\text{out}}$  denote the input and output dimensions of the layer;  $r$  is the LoRA rank with  $r \ll \min(d_{\text{in}}, d_{\text{out}})$ .

We use PrefOdds and UtilOdds throughout to track how interventions shift concept preference versus task utility on a common additive scale, with additional derivations in Appendix D.

### 3.2 Unified Dynamic Weight Formulation

We propose a unified framework that encompasses both parameter update methods and activation steering methods, by viewing them as dynamic weight update. In this formulation, both classes of methods can be expressed under a shared affine transformation view of intermediate representations (see Section 2).

Formally, the dynamic modification of the weight matrix  $\mathbf{W}$  and bias vector  $\mathbf{b}$  during inference can be written as:

$$\mathbf{h}_{i+1} = (\mathbf{W} + m_1 \Delta\mathbf{W}) \mathbf{h}_i + (\mathbf{b} + m_2 \Delta\mathbf{b}), \quad (7)$$

where  $\Delta\mathbf{W}$  and  $\Delta\mathbf{b}$  are update terms, and  $m_1, m_2$  are scalar scaling coefficients controlling their magnitudes. In other words, the original parameters are updated to  $\mathbf{W}' = \mathbf{W} + m_1 \Delta\mathbf{W}$  and  $\mathbf{b}' = \mathbf{b} + m_2 \Delta\mathbf{b}$  before computing the next-layer activation.

When a model weight is modified, the effect can be equivalently interpreted from the activation perspective, as a change to the activation at the corresponding position. In this view, diverse intervention methods are unified as adding a change term to the activation:

$$\Delta\mathbf{h} = m_1 \Delta\mathbf{W} \mathbf{h}_i + m_2 \Delta\mathbf{b}. \quad (8)$$

Under this unified view, local weight fine-tuning, LoRA, and activation steering are all specific instances, differing only in which components are updated: local weight fine-tuning modifies both  $\mathbf{W}$  and  $\mathbf{b}$ ; LoRA modifies  $\mathbf{W}$  via low-rank factors; activation steering modifies only  $\mathbf{b}$ . Table 1 summarizes their affine forms, corresponding activation update, and parameter sizes.

Notably, introducing explicit scaling coefficients extends traditional formulations and enables contin-

uous control over perturbation strength, a capability that plays a central role in our subsequent analysis.

### 3.3 Unified Dynamics Observation

**Experimental Setup.** We evaluate dynamic interventions on two types of tasks: (i) a personality tendency classification task (*Psychopathy*), and (ii) open-ended generation using *PowerSeeking* and the top 10 concept subsets from *AxBench*. We run experiments on Gemma-2-9B-IT and Qwen-2.5-7B-Instruct, and consider three intervention types: local weight updates, low-rank adaptation LoRA, and vector-based updates. We train each type using both the language modeling objective and the RePS objective. Additionally, for vector-based updates, we include a train-free method called DiffMean.

**Metrics.** For each query  $q$  with matched answers  $(A_p, A_n)$ , we compute preference and utility using the log-odds in Eqs. (5) and (6). These metrics allow us to track how preference and utility evolve as we vary the intervention scale  $m$ .

**Unified Dynamics.** Experimental results show that, under the unified perspective framework, different intervention forms exhibit remarkably consistent dynamic patterns. As shown in Figure 2, localized weight updates, low-rank adaptation (LoRA), and vector-based interventions display highly similar overall curve shapes.

For preference log-odds, all methods typically follow a three-stage pattern when plotted against the steering factor  $m$ : for small  $|m|$ , they enter a *Linear Region*, where log-odds grows approximately linearly with  $m$  (Bigelow et al., 2025); this is followed by a *Transitional Region* with a noticeable change in trend, and finally a *Convergence Region* where the curve flattens and stabilizes.

Utility log-odds, in contrast, generally peak near  $m \approx 0$ , and remain near their maximum within this narrow range. As  $|m|$  increases, utility gradually declines and eventually stabilizes.

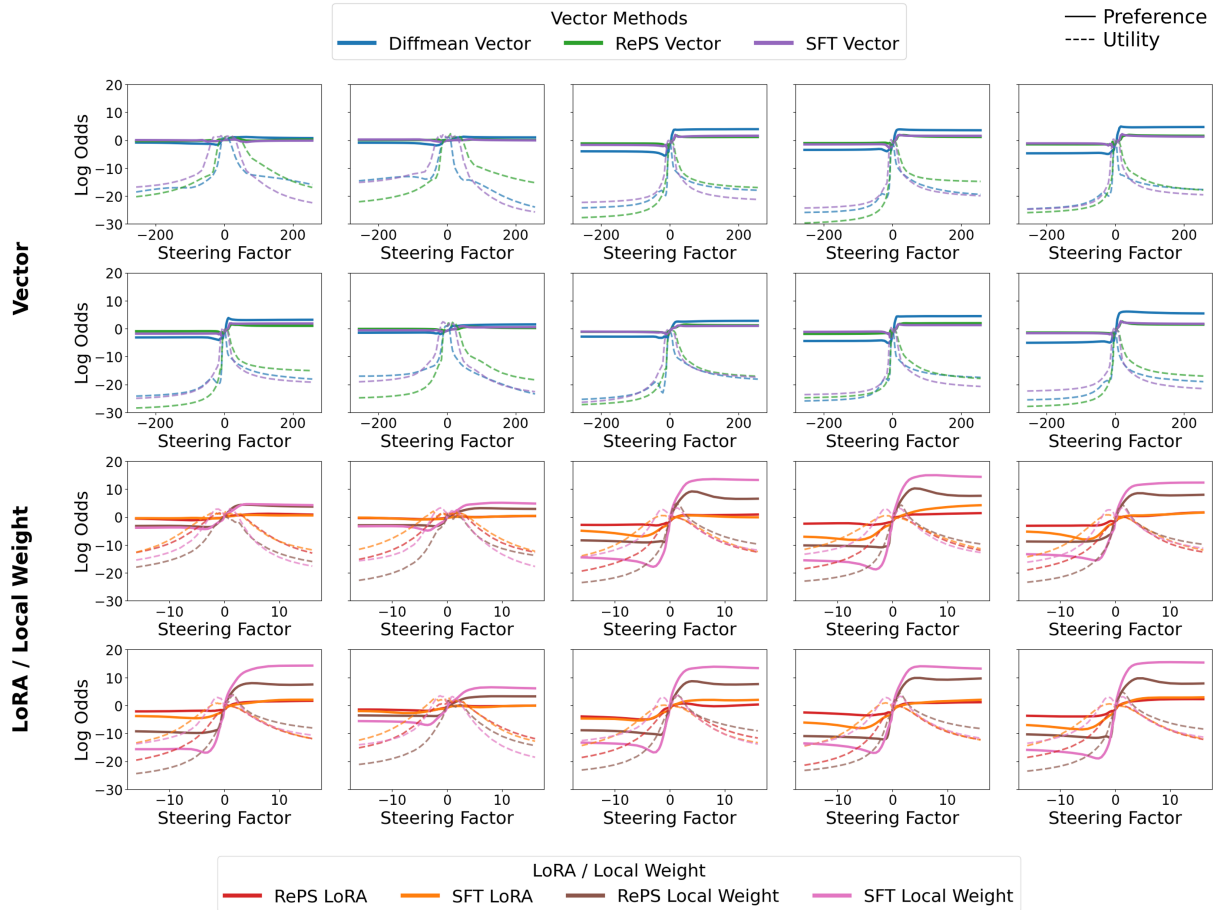


Figure 2: **Unified preference and utility dynamics under steering.** Solid lines represent preference log-odds, and dashed lines represent utility log-odds. The top panel shows steering with vector-form parameter modifications, and the bottom panel shows parametric interventions including LoRA and local weight updates. Results are shown for the Gemma-2-9B-IT model on the *AxBench* dataset, evaluated over its top 10 concept subsets. The horizontal axis corresponds to the steering factor.

325 These patterns reveal a unified steering response  
326 of preference and utility.

#### 327 4 Capability Dynamics: Mechanism 328 Analysis and Optimization

329 Motivated by the unified preference–utility dy-  
330 namics observed across intervention forms (Fig-  
331 ure 2), this section provides a mechanistic account  
332 and an empirical characterization. We take an  
333 activation-manifold perspective and introduce a  
334 simple validity-decay factor to capture the tendency  
335 for capability to degrade as steering pushes activa-  
336 tions away from the activation manifold, without  
337 committing to a specific underlying geometry. On  
338 this basis, we express preference as the combined  
339 effect of (i) steering-induced preference projection  
340 changes and (ii) validity decay, while utility is mod-  
341 eled as being dominated by the validity decay term.  
342 Finally, under this hypothesis we formalize how

343 the steering factor  $m$  shapes both preference and  
344 utility log-odds, and show via curve-fitting that the  
345 resulting forms match the observed log-odds– $m$   
346 dynamics well across settings.

#### 347 4.1 Activation Manifold Hypothesis

348 Prior work suggests that model activations of-  
349 ten concentrate on low-dimensional, manifold-like  
350 sets in representation space (Bricken et al., 2023;  
351 Wollschläger et al., 2025). Adopting this manifold  
352 perspective, we analyze additive steering as a trans-  
353 lation of hidden states along an approximately fixed  
354 direction in activation space. Intuitively, small  
355 translations may adjust model behavior in a tar-  
356 geted way, whereas large translations may push  
357 representations away from the high-density region  
358 learned during training, increasing the risk of a  
359 representation–decoder mismatch and thus degrad-  
360 ing general capability.

361 We formalize this view with two assumptions.

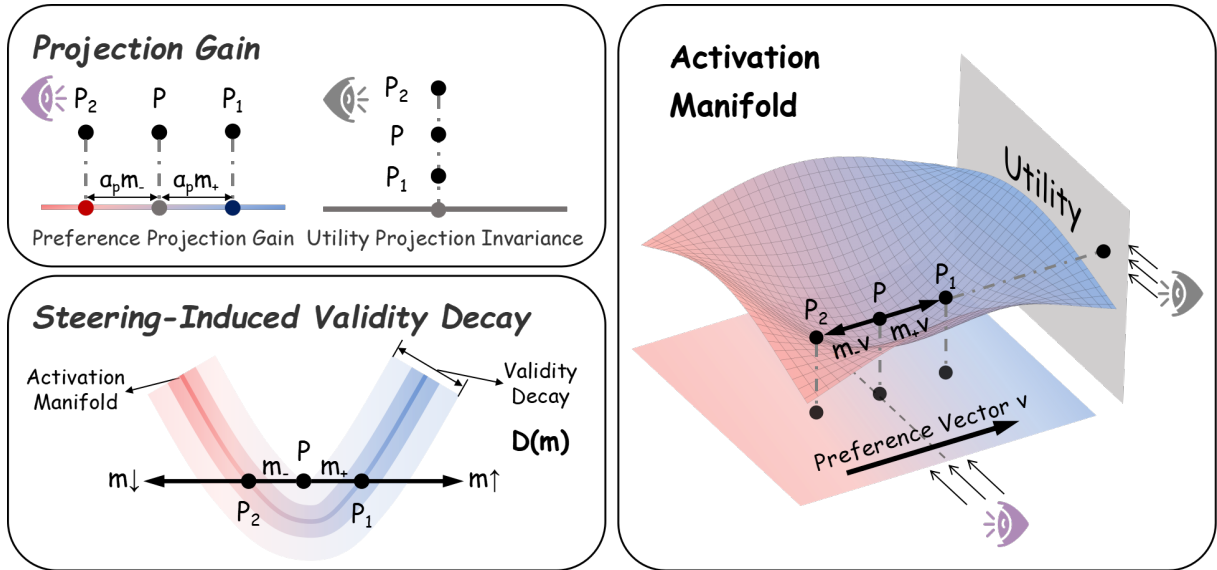


Figure 3: **Mechanism of projection gain and validity decay.** **Right:** An activation manifold view illustrating Assumption 4.1. An activation  $P$  lies on or near the manifold. Steering using preference vector  $v$  with scaling factors  $m_+$  and  $m_-$  moves  $P$  to  $P_1$  and  $P_2$ , corresponding to intersections with the manifold. **Top-left:** Projection gain. Projections onto the utility axis exhibit limited variation, whereas projections along the preference direction differ between  $P_1$  and  $P_2$ , suggesting that steering primarily influences preference-related components. **Bottom-left:** Steering-induced validity decay. As assumed in Assumption 4.2, increasing steering factor increases off-manifold deviation, leading to a monotonic decrease in validity and degraded downstream decoding.

**Assumption 4.1 (Training-Induced Activation Manifold).** Fix a layer  $l$  with hidden dimension  $d_l$ . There exists a low-dimensional set (or its neighborhood)  $\mathcal{M}_l \subset \mathbb{R}^{d_l}$  such that for inputs  $x$  drawn from a set of *stably handled* inputs  $\mathcal{X}_{\text{stable}}$ , the corresponding activation  $h_l(x)$  lies on or near  $\mathcal{M}_l$  with high probability:

$$\Pr_{x \sim \mathcal{X}_{\text{stable}}} [d(h_l(x), \mathcal{M}_l) \leq \epsilon] \geq 1 - \delta, \quad (9)$$

where  $d(\cdot, \mathcal{M}_l)$  denotes distance to  $\mathcal{M}_l$ ,  $\epsilon > 0$  is a neighborhood radius, and  $\delta \in (0, 1)$  is a small failure probability.

Assumption 4.1 asserts that pre-training induces a “typical” region of activation space where representations concentrate for stably handled inputs. We next introduce a generic notion of *representation validity*, which is high near  $\mathcal{M}_l$  and decreases as hidden states move away from it. This abstraction avoids committing to a specific geometry for  $\mathcal{M}_l$  while retaining the key implication: sufficiently off-manifold activations are more likely to be decoded unreliably by the remaining network.

**Assumption 4.2 (Steering-Induced Validity Decay).** Let  $F_{l \rightarrow L}$  denote the remainder of the model from layer  $l$  to the output logits. There exists a *validity* function  $V_l : \mathbb{R}^{d_l} \rightarrow [0, 1]$  that is monotonically non-increasing in  $d(h, \mathcal{M}_l)$ , capturing how well  $F_{l \rightarrow L}$  can stably decode an activation  $h$ .

For an additive steering intervention at layer  $l$ ,

$$\tilde{h}_l(m) = h_l + m \Delta h, \quad (10)$$

with steering direction  $\Delta h$  and steering factor  $m \in \mathbb{R}$ , define the average validity over stably handled inputs:

$$D(m) \triangleq \mathbb{E}_{x \sim \mathcal{X}_{\text{stable}}} [V_l(\tilde{h}_l(m))]. \quad (11)$$

We assume that  $D(m) \in [0, 1]$  decreases with  $|m|$  (i.e., larger interventions induce larger off-manifold shifts on average), and that the resulting capability degradation is dominated by this validity decay.

To connect Assumptions 4.1–4.2 to a concrete functional form, we view steering as moving an activation along a one-dimensional line in representation space,  $\tilde{h}_l(m) = h_l + m \Delta h$ . Under the

manifold hypothesis, degradation is governed primarily by how far this line trajectory departs from the typical region near  $\mathcal{M}_l$ , so it is natural to model  $D(m)$  as a smooth function of the (signed) distance along this line to the nearest “on-manifold” locations. In particular, as illustrated in Fig. 3, the steered trajectory may intersect the manifold neighborhood at one or more values  $\{m_i\}$  (e.g., one for  $m > 0$  and one for  $m < 0$ ). We therefore model validity as being highest near these intersection points and decaying as  $|m - m_i|$  grows.

A convenient choice that is positive, smooth, and exhibits heavy-tailed distance-based decay is the rational quadratic (RQ) form, widely used in kernel methods and Gaussian processes to model multi-scale, polynomial-rate attenuation with distance (Rasmussen, 2004). Prior research on controllability metrics has established that model steerability is often asymmetric (Miehling et al., 2025), exhibiting varying degrees of responsiveness along different directions of the same dimension. Motivated by this observation, we employ a piecewise parameterized model to quantify degradation:

$$D(m) = \begin{cases} \left(1 + \frac{(m-m_+)^2}{L_+}\right)^{-p_+} & \text{if } m \geq 0 \\ \left(1 + \frac{(m-m_-)^2}{L_-}\right)^{-p_-} & \text{if } m < 0 \end{cases} \quad (12)$$

where  $m_{\pm}$  corresponds to the signed distance from the original activation point  $P$  to an on-manifold intersection point  $P_{\pm}$  along the steering line (Fig. 3);  $L_{\pm}$  sets the characteristic scale of decay and reflects how fast the distance-to-manifold grows along the steering direction (larger when the direction is locally parallel to the manifold and smaller when it cuts across it); and  $p_{\pm}$  controls the decay rate (tail heaviness) as the trajectory moves away from the manifold neighborhood.

## 4.2 Preference Capability: Projection Gain With Decay

We study how additive steering changes a model’s preference through intermediate activations. An intervention at layer  $l$  updates the hidden state as  $\tilde{h}(m) = h + m \Delta h$ .

Prior work under LRH-style assumptions often models preference probability with a logistic form,  $P(p_p | h) = \sigma(-(\omega_p^T h + b_p))$ , where  $\omega_p$  is the preference vector. Separately, work on activation geometry suggests that after low-dimensional projection (e.g., PCA), opposite preference labels are often approximately linearly separable. Un-

der the activation-manifold view, this motivates a two-dimensional *preference plane* and a preference direction whose signed coordinate reflects *preference intensity*. Our contribution is to incorporate validity attenuation  $D(\cdot)$  (Assumption 4.2) to account for off-manifold steering.

To model this, we write the steered preference probability as

$$P(p_p | \tilde{h}(m)) = \sigma\left(-(\omega_p^T h + \alpha_p m + b_p) D_p(m)\right), \quad (13)$$

where  $\alpha_p \triangleq \omega_p^T \Delta h$  measures how much the steering direction aligns with the preference vector:  $\alpha$  is large when  $\Delta h$  is aligned with  $\omega_p$ , and  $\alpha_p = 0$  when  $\Delta h$  is orthogonal to  $\omega_p$ .

This implies the preference log-odds

$$\log \frac{P(p_p | \tilde{h}(m))}{P(p_n | \tilde{h}(m))} = (\omega_p^T h + \alpha_p m + b_p) D_p(m). \quad (14)$$

**Key implication (linear regime  $\rightarrow$  non-linear collapse).** From Eq. (14), the  $m$ -dependence enters as  $\alpha_p m D_p(m)$ . When  $|m - m^{\pm}| \ll L^{\pm}$ , Eq. (12) gives  $D_p(m) \approx 1$ , hence preference log-odds is approximately linear in  $m$  with slope  $\alpha$  (matching the near-linear regime in Bigelow et al. (2025)). As  $|m - m^{\pm}|$  grows and becomes comparable to or larger than  $L^{\pm}$ , Eq. (12) implies substantial decay in  $D_p(m)$ , so attenuation dominates and the log-odds response becomes strongly nonlinear and can collapse off-manifold.

**Fitting Form.** We fit the measured preference log-odds as a function of  $m$  with

$$\log \frac{P(p_p | \tilde{h}(m))}{P(p_n | \tilde{h}(m))} = (\alpha m + \beta) D(m) + b, \quad (15)$$

where  $\beta = \omega_p^T h + b_p$  is a per-example constant (since  $h$  is fixed for a given input), and  $b$  is an offset.

**Fit Results.** Table 2 reports the fit quality of Eq. (15), with  $R^2$  values exceeding 0.95 across most settings. These results validate the model’s ability to accurately characterize the dynamics of preference log-odds.

## 4.3 Utility Capability: Only Validity Decay

**Utility Log-odds Under Manifold-Validity Decay.** Let  $h \in \mathbb{R}^{d_l}$  denote the activation at layer

l. We quantify utility capability by the log-odds of positive vs. negative utility outcomes ( $u_p/u_n$ ). Similar to preference, we assume utility is also associated with a direction  $\omega_u$  in activation space. Under steering  $\tilde{h}(m) = h + m \Delta h$ , we model

$$\log \frac{P(u_p | \tilde{h}(m))}{P(u_n | \tilde{h}(m))} \approx -(\omega_u^\top \mathbf{h} + b_u) D_u(m), \quad (16)$$

where  $D_u(m)$  follows the manifold-validity decay in Eq. (12) and decreases with  $|m|$ . Crucially, for *preference* steering directions we typically have  $\omega_u^\top \Delta \mathbf{h} \approx 0$ , so utility is affected primarily through validity decay rather than a direct projection term.

**Fitting form.** Accordingly, we fit the measured utility log-odds with a pure decay curve:

$$\log \frac{P(u_p | \tilde{h}(m))}{P(u_n | \tilde{h}(m))} = \beta D_u(m) + b, \quad (17)$$

where  $\beta$  is the baseline log-odds and  $b$  is an offset capturing residual bias.

**Fit Results.** Table 2 reports the fit quality of Eq. (17). Uniformly high  $R^2$  values (typically  $> 0.97$ ) suggest utility variations under preference steering are well captured by the proposed formulation. Additional details are in Appendix E.

## 5 Method

### 5.1 Preference–Utility Joint Optimization

Building on the preceding mechanistic analysis, we propose a concise objective that jointly optimizes a model’s *preference* and *utility*. Our goal is to improve preference while delaying utility degradation by extending the effective linear regime of activation intervention.

**Utility Loss.** To preserve utility, we train on both the positive and negative samples for the same input using the language-modeling cross-entropy:

$$\mathcal{L}_{\text{util}} = \lambda_p \mathcal{L}_p + \lambda_n \mathcal{L}_n, \quad (18)$$

where  $\mathcal{L}_p$  and  $\mathcal{L}_n$  are the token-level cross-entropy losses on positive and negative samples, respectively, and  $\lambda_p, \lambda_n$  control their relative weight.

**Preference Loss.** By Eq. (5), the loss gap  $\mathcal{L}_n - \mathcal{L}_p$  is exactly the preference log-odds. We therefore maximize this gap via a hinge-style margin loss:

$$\mathcal{L}_{\text{pref}} = \gamma \cdot \sigma(\theta - (\mathcal{L}_n - \mathcal{L}_p)), \quad (19)$$

where  $\sigma(\cdot)$  is ReLU and  $\theta$  is a margin threshold, and  $\gamma$  trades off preference improvement against utility preservation.

Type	Method	Preference $R^2 \uparrow$				Utility $R^2 \uparrow$			
		PSY	PWR	AXB	Avg	PSY	PWR	AXB	Avg
<b>Gemma-2-9B-IT</b>									
Weight	SFT	0.97	0.98	0.99	0.98	0.98	0.99	0.98	0.98
	RePS	0.99	0.99	0.99	0.99	0.96	0.99	0.99	0.98
LoRA	SFT	0.92	0.99	0.98	0.96	0.98	0.99	0.99	0.99
	RePS	0.83	0.99	0.99	0.94	0.99	0.99	0.99	0.99
Vector	DiffMean	0.97	0.99	0.99	0.98	0.97	0.99	0.98	0.98
	SFT	0.93	0.97	0.98	0.96	0.99	0.99	0.99	0.99
	RePS	0.99	0.98	0.95	0.97	0.99	0.99	0.99	0.99
<b>Qwen-2.5-7B-IT</b>									
Weight	SFT	0.99	0.99	0.99	0.99	0.99	0.99	0.98	0.99
	RePS	0.99	0.99	0.99	0.99	0.99	0.99	0.97	0.98
LoRA	SFT	0.97	0.99	0.99	0.98	0.99	0.99	0.99	0.99
	RePS	0.94	0.99	0.99	0.97	0.98	0.96	0.99	0.98
Vector	DiffMean	0.98	0.95	0.98	0.97	0.99	0.99	0.97	0.98
	SFT	0.93	0.97	0.98	0.96	0.98	0.99	0.99	0.99
	RePS	0.97	0.98	0.93	0.96	0.99	0.98	0.99	0.99

Table 2: **Curve fitting performance.** Results on Psychopathy (PSY), PowerSeeking (PWR), and AXBench (AXB). We report  $R^2$  (higher is better), measuring alignment between theoretical curves and empirical data. Color intensity indicates  $R^2$  values. Consistently dark shading shows high fidelity across settings ( $R^2 > 0.95$ ).

**Final Objective.** We combine the two components as

$$\mathcal{L} = \mathcal{L}_{\text{util}} + \mathcal{L}_{\text{pref}}. \quad (20)$$

### 5.2 Experiment Results.

We evaluate the proposed preference-utility joint optimization method under three intervention forms: local weight update, low-rank adaptation (LoRA), and activation vector steering. As shown in Table 3 at Appendix B, our approach consistently achieves higher scores compared with baseline methods across all three intervention types. These results demonstrate the robustness and generality of the proposed optimization strategy.

## 6 Conclusion

We propose a unified dynamic weight update framework that incorporates parameter updates, LoRA, and activation interventions, revealing a consistent preference–utility decay pattern in the log-odds space. Building on this mechanistic insight, we design a joint optimization method that consistently improves preference while mitigating utility degradation across diverse intervention forms, demonstrating versatility and robustness.

## 521 Limitations

522 While our unified dynamic weight update frame-  
523 work provides a coherent perspective on LLM  
524 control and enables predictable preference–utility  
525 trade-offs, several limitations remain. First, our  
526 analysis assumes that model representations lie  
527 near a well-structured activation manifold, which  
528 may not hold for extremely large or highly di-  
529 verse models, potentially reducing the accuracy  
530 of our quantitative predictions. Second, our exper-  
531 iments focus primarily on attribute-level control  
532 (e.g., sentiment, style), leaving the applicability to  
533 complex multi-turn reasoning or safety-critical con-  
534 tent largely unexplored. Third, while our proposed  
535 training objective mitigates the utility–preference  
536 trade-off, it does not guarantee complete avoid-  
537 ance of undesirable side effects such as subtle in-  
538 struction violations or context drift under extreme  
539 control strengths. Finally, our study evaluates con-  
540 trol under pre-defined intervention multipliers, and  
541 generalization to adaptive or dynamically varying  
542 control signals requires further investigation.

## 543 Ethics Statement

544 Controlled LLM generation carries inherent ethical  
545 considerations. While our framework aims to im-  
546 prove controllability and preserve task validity, it  
547 could potentially be misused to manipulate user per-  
548 ception, amplify biased viewpoints, or generate per-  
549 suasive yet misleading content. Our experiments  
550 are conducted on standard benchmark datasets and  
551 do not involve sensitive personal information. We  
552 emphasize that the proposed methods should be  
553 deployed with human oversight, adherence to fair-  
554 ness guidelines, and robust monitoring to prevent  
555 harm. By explicitly modeling preference–utility  
556 trade-offs, we aim to make LLM interventions  
557 more interpretable and safer, but responsible use  
558 depends on context-aware implementation and  
559 alignment with societal norms.

## 560 References

561 Amrita Bhattacharjee, Shaona Ghosh, Traian Rebedea,  
562 and Christopher Parisien. 2024. [Towards inference-](#)  
563 [time category-wise safety steering for large language](#)  
564 [models](#). *CoRR*, abs/2410.01174.

565 Eric J. Bigelow, Daniel Wurgaft, YingQiao Wang,  
566 Noah D. Goodman, Tomer D. Ullman, Hidenori  
567 Tanaka, and Ekdeep Singh Lubana. 2025. [Belief](#)  
568 [dynamics reveal the dual nature of in-context learn-](#)  
569 [ing and activation steering](#). *CoRR*, abs/2511.00617.

Trenton Bricken, Adly Templeton, Joshua Batson, 570  
Brian Chen, Adam Jermy, Tom Conerly, Nick 571  
Turner, Cem Anil, Carson Denison, Amanda Askill, 572  
Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas 573  
Schiefer, Tim Maxwell, Nicholas Joseph, Zac 574  
Hatfield-Dodds, Alex Tamkin, Karina Nguyen, 575  
Brayden McLean, Josiah E Burke, Tristan Hume, 576  
Shan Carter, Tom Henighan, and Christopher 577  
Olah. 2023. [Towards monosemanticity: Decom-](#)  
578 [posing language models with dictionary learning](#).  
579 *Transformer Circuits Thread*. [https://transformer-](https://transformer-circuits.pub/2023/monosemantic-features/index.html)  
580  [582  \[584  \\[586  \\\[588  \\\\[590\\\\]\\\\(https://transformer-circuits.pub/2023/monosemantic-</a><br/>589 <a href=\\\\)\\\]\\\(https://transformer-circuits.pub/2023/monosemantic-</a><br/>587 <a href=\\\)\\]\\(https://transformer-circuits.pub/2023/monosemantic-</a><br/>585 <a href=\\)\]\(https://transformer-circuits.pub/2023/monosemantic-</a><br/>583 <a href=\)](https://transformer-circuits.pub/2023/monosemantic-</a><br/>581 <a href=)

Yuanpu Cao, Tianrong Zhang, Bochuan Cao, Ziyi Yin, 583  
Lu Lin, Fenglong Ma, and Jinghui Chen. 2024. [Perso-](#)  
584 [nalized steering of large language models: Versa-](#)  
585 [tile steering vectors through bi-directional preference](#)  
586 [optimization](#). In *Advances in Neural Information*  
587 *Processing Systems 38: Annual Conference on Neu-*  
588 *ral Information Processing Systems 2024, NeurIPS*  
589 *2024, Vancouver, BC, Canada, December 10 - 15,*  
590 *2024*. 591

Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zong- 592  
han Yang, Yusheng Su, Shengding Hu, Yulin Chen, 593  
Chi-Min Chan, Weize Chen, Jing Yi, Weilin Zhao, 594  
Xiaozhi Wang, Zhiyuan Liu, Hai-Tao Zheng, Jianfei 595  
Chen, Yang Liu, Jie Tang, Juanzi Li, and Maosong 596  
Sun. 2022. [Delta tuning: A comprehensive study of](#)  
597 [parameter efficient methods for pre-trained language](#)  
598 [models](#). *CoRR*, abs/2203.06904. 599

Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, 600  
Zonghan Yang, Yusheng Su, Shengding Hu, Yulin 601  
Chen, Chi-Min Chan, Weize Chen, et al. 2023. 602  
Parameter-efficient fine-tuning of large-scale pre- 603  
trained language models. *Nature machine intelli-*  
604 *gence*, 5(3):220–235. 605

Mor Geva, Roei Schuster, Jonathan Berant, and Omer 606  
Levy. 2021. [Transformer feed-forward layers are key-](#)  
607 [value memories](#). In *Proceedings of the 2021 Confer-*  
608 *ence on Empirical Methods in Natural Language Pro-*  
609 *cessing, EMNLP 2021, Virtual Event / Punta Cana,*  
610 *Dominican Republic, 7-11 November, 2021*, pages  
611 5484–5495. Association for Computational Linguis-  
612 tics. 613

Chi Han, Jialiang Xu, Manling Li, Yi Fung, Chenkai 614  
Sun, Nan Jiang, Tarek Abdelzaher, and Heng Ji. 2024. 615  
Word embeddings are steers for language models. 616  
In *Proceedings of the 62nd Annual Meeting of the*  
617 *Association for Computational Linguistics (Volume*  
618 *1: Long Papers)*, pages 16410–16430. 619

Peixuan Han, Cheng Qian, Xiushi Chen, Yuji Zhang, 620  
Denghui Zhang, and Heng Ji. 2025. Internal activa- 621  
tion as the polar star for steering unsafe llm behavior. 622  
*arXiv preprint arXiv:2502.01042*. 623

Soufiane Hayou, Nikhil Ghosh, and Bin Yu. 2024. 624  
[Lora+: Efficient low rank adaptation of large models](#). 625  
In *Forty-first International Conference on Machine*  
626 *Learning, ICML 2024, Vienna, Austria, July 21-27,*  
627 *2024*. OpenReview.net. 628



742	Curt Tigges, Curt Tigges, Oskar Hollinsworth, Curt	<i>Short Papers</i> ), <i>ACL 2022, Dublin, Ireland, May 22-</i>	799
743	Tigges, Atticus Geiger, Atticus Geiger, Oskar	27, 2022, pages 1–9. Association for Computational	800
744	Hollinsworth, Neel Nanda, Neel Nanda, Atticus	Linguistics.	801
745	Geiger, and Neel Nanda. 2023. <a href="#">Linear repre-</a>		
746	<a href="#">sentations of sentiment in large language models.</a>	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,	802
747	<a href="http://arxiv.org/abs/2310.15154">http://arxiv.org/abs/2310.15154</a> .	Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen	803
		Zhang, Junjie Zhang, Zican Dong, et al. 2023. A	804
748	Alexander Matt Turner, Lisa Thiergart, Gavin Leech,	survey of large language models. <i>arXiv preprint</i>	805
749	David Udell, Juan J Vazquez, Ulisse Mini, and Monte	<i>arXiv:2303.18223</i> , 1(2).	806
750	MacDiarmid. 2023. Activation addition: Steering		
751	language models without optimization. <i>arXiv e-</i>	Bojia Zi, Xianbiao Qi, Lingzhi Wang, Jianan Wang,	807
752	<i>prints</i> , pages arXiv–2308.	Kam-Fai Wong, and Lei Zhang. 2023. <a href="#">Delta-lora:</a>	808
		<a href="#">Fine-tuning high-rank parameters with the delta of</a>	809
753	Teun van der Weij, Massimo Poesio, and Nandi Schoots.	<a href="#">low-rank matrices.</a> <i>CoRR</i> , abs/2309.02411.	810
754	2024. <a href="#">Extending activation steering to broad skills</a>		
755	<a href="#">and multiple behaviours.</a> <i>CoRR</i> , abs/2403.05767.		
756	Mengru Wang, Ziwon Xu, Shengyu Mao, Shumin Deng,		
757	Zhaopeng Tu, Huajun Chen, and Ningyu Zhang.		
758	2025. <a href="#">Beyond prompt engineering: Robust behavior</a>		
759	<a href="#">control in llms via steering target atoms.</a> In <i>Proceed-</i>		
760	<i>ings of the 63rd Annual Meeting of the Association</i>		
761	<i>for Computational Linguistics (Volume 1: Long Pa-</i>		
762	<i>pers)</i> , <i>ACL 2025, Vienna, Austria, July 27 - August 1,</i>		
763	2025, pages 23381–23399. Association for Computa-		
764	tional Linguistics.		
765	Tom Wollschläger, Jannes Elstner, Simon Geisler, Vin-		
766	cent Cohen-Addad, Stephan Günemann, and Jo-		
767	hannes Gasteiger. 2025. <a href="#">The geometry of refusal</a>		
768	<a href="#">in large language models: Concept cones and repre-</a>		
769	<a href="#">sentational independence.</a> In <i>Forty-second Interna-</i>		
770	<i>tional Conference on Machine Learning, ICML 2025,</i>		
771	<i>Vancouver, BC, Canada, July 13-19, 2025.</i> OpenRe-		
772	view.net.		
773	Zhengxuan Wu, Aryaman Arora, Atticus Geiger, Zheng		
774	Wang, Jing Huang, Dan Jurafsky, Christopher D Man-		
775	ning, and Christopher Potts. 2025a. <a href="#">Axbench: Steering</a>		
776	<a href="#">llms? even simple baselines outperform sparse</a>		
777	<a href="#">autoencoders.</a> <i>arXiv preprint arXiv:2501.17148</i> .		
778	Zhengxuan Wu, Qinan Yu, Aryaman Arora, Christo-		
779	pher D. Manning, and Christopher Potts. 2025b. <a href="#">Im-</a>		
780	<a href="#">proved representation steering for language models.</a>		
781	<i>CoRR</i> , abs/2505.20809.		
782	Zhenda Xie, Yixuan Wei, Huanqi Cao, Chenggang		
783	Zhao, Chengqi Deng, Jiashi Li, Damai Dai, Huazuo		
784	Gao, Jiang Chang, Liang Zhao, Shangyan Zhou,		
785	Zhean Xu, Zhengyan Zhang, Wangding Zeng,		
786	Shengding Hu, Yuqing Wang, Jingyang Yuan, Lean		
787	Wang, and Wenfeng Liang. 2025. <a href="#">mhc: Manifold-</a>		
788	<a href="#">constrained hyper-connections.</a> <i>arXiv preprint</i>		
789	<a href="https://arxiv.org/pdf/2512.24880">https://arxiv.org/pdf/2512.24880</a> .		
790	Wanli Yang, Fei Sun, Rui Tang, Hongyu Zang, Du Su,		
791	Qi Cao, Jingang Wang, Huawei Shen, and Xueqi		
792	Cheng. 2025. <a href="#">Fine-tuning done right in model edit-</a>		
793	<a href="#">ing.</a> <i>CoRR</i> , abs/2509.22072.		
794	Elad Ben Zaken, Yoav Goldberg, and Shauli Ravfogel.		
795	2022. <a href="#">Bitfit: Simple parameter-efficient fine-tuning</a>		
796	<a href="#">for transformer-based masked language-models.</a> In		
797	<i>Proceedings of the 60th Annual Meeting of the As-</i>		
798	<i>sociation for Computational Linguistics (Volume 2:</i>		

## A Related Work

**Mechanism.** Most activation steering methods assume linear structure in activation space, controlling concepts by adding scaled direction vectors to hidden states (Mikolov et al., 2013; Pennington et al., 2014; Nanda et al., 2023; Tigges et al., 2023; Park et al., 2024). Building on this view, Bigelow et al. (2025) show that steering yields an approximately linear trend in posterior odds, but mainly in the small-scale regime. Recent studies further report non-monotonic or adverse effects under stronger steering, challenging a naive global linearity assumption (Bricken et al., 2023; Wollschläger et al., 2025). Meanwhile, representation-manifold work provides a complementary geometric lens for understanding steering and its limitations (Modell et al., 2025; Li and He, 2025; Xie et al., 2025).

**Activation Steering.** Activation steering controls the behavior of large language models (LLMs) by intervening in hidden states during forward propagation, where steering vectors are constructed to manipulate single attributes as well as more complex behavioral targets (Turner et al., 2023; Rimsky et al., 2024; van der Weij et al., 2024; Rahn et al., 2024; Scalena et al., 2024; Tan et al., 2024; Bhattacharjee et al., 2024; Postmus and Abreu, 2024; Konen et al., 2024; Hazra et al., 2024; Han et al., 2025; Jiang et al., 2025). However, recent studies have shown that the coarse-grained nature of activation steering can lead to degradation in model utility (Wang et al., 2025; Wu et al., 2025a). Cao et al. (2024); Wu et al. (2025b) introduce explicit preference learning objectives to optimize activation steering, enabling more precise control.

**Parameter-Efficient Fine-Tuning.** Parameter-Efficient Fine-Tuning (PEFT), including adapters, prompt tuning, and low-rank adaptation (LoRA), shows that effective adaptation of large language models does not require updating the full set of parameters. LoRA achieves performance comparable to full fine-tuning, suggesting that adaptation is driven by structured weight deltas rather than full parameter updates, with low-rank residuals serving as localized parameter updates (Hu et al., 2022a; Zi et al., 2023; Hayou et al., 2024; Kopiczko et al., 2024). Localized parameter update methods further indicate that knowledge and behavior in LLMs are highly localized, as factual associations can be modified by intervening on a small number of parameters in specific layers (Za-

ken et al., 2022). Ding et al. (2022) confirms that updating only a subset of parameters suffices for task adaptation with minimal impact on pretrained capabilities. Geva et al. (2021); Yang et al. (2025) propose selectively updating only the MLP layers enables stable knowledge editing and largely preserves the model’s original capabilities.

## B Experiment Details

We evaluate steering methods on three benchmarks that capture different aspects of controlled generation. The experimental results are evaluated using GPT-41-mini. Table 3 compares different intervention forms, including weight updates, LoRA, and activation vector steering, across two base models. Across all control forms, our method consistently achieves strong or best performance on concept metrics while maintaining competitive harmonic scores. Notably, the advantages of our approach are most pronounced under LoRA and activation vector interventions. In these settings, our method consistently improves concept strength over both SFT and REPS baselines, and achieves the highest or near-highest harmonic scores on AXBENCH, indicating a more favorable balance between preference enforcement and utility preservation. In contrast, gains under full weight updates are comparatively smaller, likely due to the higher expressive capacity of direct parameter modification, but our method remains competitive without sacrificing utility. Overall, these results demonstrate the generality and robustness of the proposed optimization across diverse steering mechanisms.

## C List of Mathematical Symbols

The table 4 below lists the important symbols used in this paper.

## D Derivations and Implementation Details for Log-Odds

This appendix derives the loss-based forms of Eqs. (5)–(6) from the preference–utility independence assumption, and states how we compute the required sequence losses.

Model	Intervention Form	Method	Psychopathy	PowerSeeking	AXBENCH		
			Acc $\uparrow$	Concept $\uparrow$	Concept $\uparrow$	Harmonic $\uparrow$	
Gemma-2-9B-IT	Vanilla	Vanilla	50.00	1.87	0.4750	0.4950	
		Weight	SFT	<b>100.00</b>	3.50	1.6625	1.4538
			REPS	<b>100.00</b>	3.39	<u>1.7750</u>	<b>1.6362</b>
	Ours		<b>100.00</b>	<b>3.59</b>	<b>1.8500</b>	<u>1.6225</u>	
	LoRA	SFT	<b>100.00</b>	3.41	<u>1.7625</u>	1.5188	
		REPS	<u>99.00</u>	3.44	1.7375	<b>1.6525</b>	
		Ours	<b>100.00</b>	<b>3.56</b>	<b>1.7750</b>	<u>1.6412</u>	
	Vector	DiffMean (No train)	53.00	2.95	1.1625	1.0550	
		SFT	<u>97.00</u>	3.30	<u>1.7000</u>	1.4487	
		REPS	<u>98.00</u>	<u>3.61</u>	<u>1.7000</u>	<u>1.5550</u>	
		Ours	<b>99.00</b>	<b>3.62</b>	<b>1.8500</b>	<b>1.6475</b>	
	Qwen-2.5-7B-IT	Vanilla	Vanilla	50.00	2.24	0.4500	0.4713
Weight			SFT	<u>97.00</u>	<u>3.53</u>	1.5375	1.1287
			REPS	96.00	3.24	<b>1.7750</b>	<b>1.4125</b>
		Ours	<b>98.00</b>	<b>3.66</b>	<u>1.5750</u>	<u>1.3350</u>	
LoRA		SFT	<u>99.00</u>	3.05	1.4875	1.3175	
		REPS	<b>100.00</b>	<u>3.34</u>	1.4875	1.4013	
		Ours	<b>100.00</b>	<b>3.59</b>	<b>1.7375</b>	<b>1.6362</b>	
Vector		DiffMean (No train)	55.00	3.17	0.9500	0.9950	
		SFT	<u>97.00</u>	3.58	<u>1.5750</u>	<u>1.5800</u>	
		REPS	88.00	<u>3.63</u>	<u>1.7375</u>	<u>1.6412</u>	
		Ours	<b>98.00</b>	<b>3.65</b>	<b>1.8125</b>	<b>1.6500</b>	

Table 3: Comparison of different steering methods under three control forms. All methods perform inference-time interventions on hidden representations. DiffMean is a training-free baseline constructed by mean-difference of hidden representations. Best and second-best results are highlighted within each model and intervention form.

### D.1 From preference-utility independence to log-odds

For a query  $q$  and a polarity pair  $(A_p, A_n)$ , we assume

$$\begin{aligned} P(A_p | q) &= P(u | q) P(p_p | q), \\ P(A_n | q) &= P(u | q) P(p_n | q), \end{aligned} \quad (21)$$

with  $P(p_p | q) + P(p_n | q) = 1$ .

**Preference log-odds.** Taking the ratio of (21) cancels  $P(u | q)$ :

$$\frac{P(A_p | q)}{P(A_n | q)} = \frac{P(p_p | q)}{P(p_n | q)}. \quad (22)$$

Applying  $\log(\cdot)$  gives

$$\text{PrefOdds}(q) \triangleq \log \frac{P(p_p | q)}{P(p_n | q)} = \log \frac{P(A_p | q)}{P(A_n | q)}. \quad (23)$$

Using the loss definition  $\mathcal{L} \triangleq -\log P(A | q)$ , we have  $P(A | q) = e^{-\mathcal{L}}$ , and thus

$$\text{PrefOdds}(q) = \log \frac{e^{-\mathcal{L}_p}}{e^{-\mathcal{L}_n}} = \mathcal{L}_n - \mathcal{L}_p, \quad (24)$$

which matches Eq. (5).

**Utility probability and log-odds.** Summing (21) and using  $P(p_p | q) + P(p_n | q) = 1$  yields

$$\begin{aligned} P(A_p | q) + P(A_n | q) &= P(u | q) \\ &\quad \left( P(p_p | q) + P(p_n | q) \right) \\ &= P(u | q). \end{aligned} \quad (25)$$

Therefore,

$$\begin{aligned} \text{UtilOdds}(q) &\triangleq \log \frac{P(u | q)}{1 - P(u | q)} \\ &= \log \frac{P(A_p | q) + P(A_n | q)}{1 - P(A_p | q) - P(A_n | q)}. \end{aligned} \quad (26)$$

Substituting  $P(A | q) = e^{-\mathcal{L}(A|q)}$  gives the loss form

$$\text{UtilOdds}(q) = \log \frac{e^{-\mathcal{L}_p} + e^{-\mathcal{L}_n}}{1 - e^{-\mathcal{L}_p} - e^{-\mathcal{L}_n}}, \quad (27)$$

which matches Eq. (6). Note that since  $(A_p, A_n)$  are only two candidate continuations, we typically have  $P(A_p | q) + P(A_n | q) < 1$ .

## 932 D.2 Computing sequence losses

933 Let  $A = (y_1, \dots, y_T)$  be a completion (excluding  
934 the query/prompt tokens). We compute the se-  
935 quence negative log-likelihood (cross-entropy loss)  
936 under teacher forcing:

$$937 \mathcal{L}(A | q) \triangleq -\log P(A | q) \\ 938 = -\sum_{t=1}^T \log P(y_t | q, y_{<t}). \quad (28)$$

939 We then set  $\mathcal{L}_p \triangleq \mathcal{L}(A_p | q)$  and  $\mathcal{L}_n \triangleq \mathcal{L}(A_n | q)$   
940 and plug them into Eqs. (24) and (27).

941 **Length normalization (optional).** When  $A_p$  and  
942  $A_n$  have different lengths, we optionally use the  
943 mean loss  $\bar{\mathcal{L}}(A | q) \triangleq \mathcal{L}(A | q)/T$  in place of  
944  $\mathcal{L}(A | q)$  to reduce length effects. In that case, the  
945 corresponding quantities use  $e^{-\bar{\mathcal{L}}}$  instead of  $e^{-\mathcal{L}}$ .

## 946 E Fitting Experiment Details

### 947 E.1 Fitting Results on Test Set

948 To further validate our theoretical model, we per-  
949 formed parameterized fitting on the test set using  
950 the SLSQP algorithm, strictly enforcing continuity  
951 between positive and negative segments at the ori-  
952 gin. As shown in Table 5, the direct fitting yielded  
953 high goodness-of-fit values ( $R^2 > 0.95$ ) for most  
954 methods. This confirms that the steering effect  
955 follows a deterministic trajectory predicted by our  
956 theory rather than random perturbations, thereby  
957 validating the proposed interaction mechanism.

### 958 E.2 Analysis of Generalization Ability

959 Following the validation of our theoretical mech-  
960 anism, we conducted train-to-test transfer exper-  
961 iments to evaluate the extent to which different  
962 methods decouple "concepts" from specific inputs.  
963 Theoretical curve parameters were derived solely  
964 from training data and applied directly to the test  
965 set for prediction (Table 6).

966 **Robust Generalization** Most activation-based  
967 methods maintained high prediction accuracy on  
968 the test set, indicating they successfully extract gen-  
969 eral, input-agnostic conceptual representations that  
970 adhere to universal theoretical laws across samples.

Symbol	Description
<b>Unified Analysis Framework</b>	
$m$	The steering scalar coefficient.
$\text{PrefOdds}(q)$	Preference log-odds ( $\mathcal{L}_n - \mathcal{L}_p$ ). (5)
$\text{UtilOdds}(q)$	Utility log-odds. (6)
$P(u q)$	Latent utility probability.
$P(p_{\pm} q)$	Latent preference probability.
$P(\bullet h)$	Equivalent to $P(\bullet q)$ as the weights remain unchanged.
$\mathcal{L}_p, \mathcal{L}_n$	Cross-entropy losses correspond- ing to $A_p$ and $A_n$ .
<b>Mechanistic Manifold Model</b>	
$\mathcal{M}_l$	The activation manifold of stably handled inputs at layer $l$ .
$D(m)$	Average validity decay function.
$m_{\pm}$	Distance from P to P $\pm$ along the steering line. 3
$L_{\pm}$	Characteristic scale of decay.
$p_{\pm}$	Asymptotic decay rate.
<b>Joint Optimization</b>	
$\mathcal{L}_{util}$	Utility loss component (standard language modeling loss).
$\mathcal{L}_{pref}$	Preference loss component (margin-based hinge loss).
$\gamma$	Trade-off coefficient balancing preference improvement and util- ity preservation.

Table 4: **Notations for Key Concepts.** A summary of the specialized symbols introduced for the unified analysis, mechanistic modeling, and optimization objective.

Type	Method	Preference $R^2 \uparrow$				Utility $R^2 \uparrow$			
		PSY	PWR	AXB	Avg	PSY	PWR	AXB	Avg
<i>Gemma-2-9B-IT</i>									
Weight	SFT	0.96	0.96	0.99	0.97	0.98	0.93	0.99	0.97
	RePS	0.95	0.98	0.95	0.96	0.98	0.93	0.99	0.97
LoRA	SFT	0.99	0.99	0.98	0.99	0.98	0.98	0.99	0.98
	RePS	0.99	0.99	0.98	0.99	0.99	0.99	0.99	0.99
Vector	DiffMean	0.89	0.99	0.99	0.96	0.94	0.99	0.98	0.97
	SFT	0.90	0.97	0.97	0.95	0.98	0.99	0.99	0.99
	RePS	0.96	0.98	0.96	0.97	0.96	0.99	0.99	0.98
<i>Qwen-2.5-7B-IT</i>									
Weight	SFT	0.99	0.82	0.99	0.93	0.99	0.99	0.95	0.98
	RePS	0.99	0.89	0.97	0.95	0.98	0.99	0.90	0.96
LoRA	SFT	0.70	0.95	0.98	0.88	0.99	0.99	0.99	0.99
	RePS	0.88	0.95	0.95	0.93	0.98	0.99	0.98	0.98
Vector	DiffMean	0.99	0.99	0.98	0.99	0.97	0.94	0.98	0.96
	SFT	0.99	0.99	0.97	0.98	0.97	0.95	0.99	0.97
	RePS	0.99	0.98	0.93	0.97	0.96	0.96	0.98	0.97

Table 5: **Performance comparison of curve fitting quality on test sets.** We evaluate the models on three datasets: Psychopathy (PSY), PowerSeeking (PWR), and AXBench (AXB).

Type	Method	Preference $R^2 \uparrow$				Utility $R^2 \uparrow$			
		PSY	PWR	AXB	Avg	PSY	PWR	AXB	Avg
<i>Gemma-2-9B-IT</i>									
Weight	SFT	0.96	0.85	-4.25	-0.81	0.98	0.98	0.61	0.86
	RePS	0.99	0.98	-1.16	0.27	0.96	0.98	0.73	0.89
LoRA	SFT	0.92	0.99	-0.56	0.45	0.98	0.99	0.96	0.98
	RePS	0.83	0.99	0.74	0.85	0.98	0.99	0.97	0.98
Vector	DiffMean	-0.14	0.99	0.75	0.53	0.97	0.99	0.97	0.98
	SFT	0.90	0.91	0.74	0.85	0.98	0.99	0.99	0.99
	RePS	0.98	0.89	0.65	0.84	0.99	0.99	0.99	0.99
<i>Qwen-2.5-7B-IT</i>									
Weight	SFT	0.99	-0.32	-12.03	-3.79	0.99	-1.33	-3.07	-1.14
	RePS	0.96	0.98	-3.82	-0.63	0.99	0.42	-1.15	0.09
LoRA	SFT	0.97	0.98	-0.40	0.52	0.99	0.99	0.95	0.98
	RePS	0.94	0.99	-0.13	0.60	0.98	0.96	0.96	0.97
Vector	DiffMean	0.86	0.94	0.80	0.87	0.37	0.99	0.97	0.78
	SFT	0.67	0.92	0.71	0.77	0.96	0.99	0.99	0.98
	RePS	0.97	0.93	0.74	0.88	0.99	0.98	0.98	0.98

Table 6: **Generalization ability of curve fitting.** The table reports the  $R^2$  scores where the curves are fitted on the training set and evaluated on the test set across three datasets: Psychopathy (PSY), PowerSeeking (PWR), and AXBench (AXB). Negative values imply that the fitted curves do not generalize well to unseen data.