
Angular Steering: Behavior Control via Rotation in Activation Space

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Controlling specific behaviors in large language models while preserving their
2 general capabilities is a central challenge for safe and reliable artificial intelligence
3 (AI) deployment. Current steering methods, such as vector addition and direc-
4 tional ablation, are constrained within a two-dimensional subspace defined by the
5 activation and feature direction, making them sensitive to chosen parameters and
6 potentially affecting unrelated features due to unintended interactions in activation
7 space. We introduce Angular Steering, a novel and flexible method for behavior
8 modulation that operates by rotating activations within a fixed two-dimensional
9 subspace. By formulating steering as a geometric rotation toward or away from
10 a target behavior direction, Angular Steering provides continuous, fine-grained
11 control over behaviors such as refusal and compliance. We demonstrate this method
12 using refusal steering as a use case. Additionally, we propose Adaptive Angular
13 Steering, a selective variant that rotates only activations aligned with the target
14 feature, further enhancing stability and coherence. Angular Steering generalizes
15 existing addition and orthogonalization techniques under a unified geometric rota-
16 tion framework, simplifying parameter selection and maintaining model stability
17 across a broader range of adjustments. Experiments across multiple model families
18 and sizes show that Angular Steering achieves robust behavioral control while
19 maintaining general language modeling performance, underscoring its flexibility,
20 generalization, and robustness compared to prior approaches.

21 1 Introduction

22 Large language models (LLMs) have become remarkably capable, yet steering their behavior towards
23 desired responses remains a challenge. On one hand, we want the model to follow certain guidelines
24 or exhibit particular traits, e.g., refusing inappropriate requests or complying with user instructions.
25 On the other hand, aggressive tuning of the models behavior can degrade its original performance,
26 causing losses in fluency or actuality [45, 47].

27 Activation steering, which manipulates internal representations of language models at inference time,
28 has emerged as a compelling alternative to retraining for behavior control [47, 54, 35]. Techniques
29 such as activation addition [47, 35] and direction orthogonalization [1, 54] have demonstrated the
30 capacity to steer models toward or away from specific behaviors. However, these methods offer limited
31 granularity. For instance, orthogonalization removes the feature entirely by projecting activations
32 onto the orthogonal subspace, leaving no room for partial suppression. Moreover, activation addition
33 requires careful tuning of the coefficient to avoid instability; improper values can lead to degraded
34 fluency or incoherent outputs [38, 43, 48, 39]. While conditional methods improve context-sensitivity,
35 they often retain the underlying manipulation mechanism [50, 18, 20].

36 **Contribution.** We propose *Angular Steering*, a method that reformulates activation editing as a
37 geometric rotation within a fixed 2D subspace of the models activation space. By identifying a
38 behavior-associated direction, e.g., refusal, and a complementary axis that jointly define a steering

Submitted to 39th Conference on Neural Information Processing Systems (NeurIPS 2025). Do not distribute.

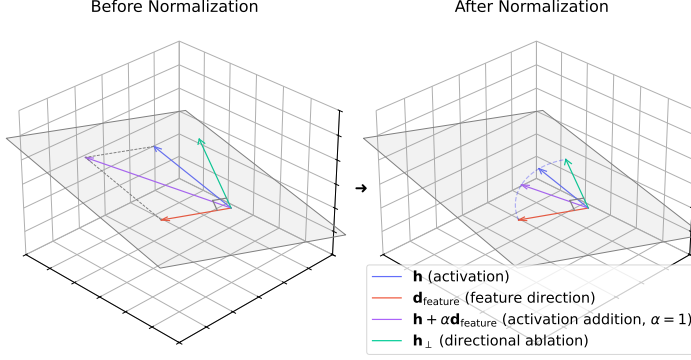


Figure 1: Geometric interpretation of activation steering. Left: Before normalization, the original activation vector \mathbf{h} , the feature direction \mathbf{d}_{feat} , the ablation vector \mathbf{h}_{\perp} , and addition $\mathbf{h} + \alpha \mathbf{d}_{\text{feat}}$ are shown. Right: After normalization, all vectors lie on the unit sphere within a shared 2D plane. The dashed arc shows \mathbf{h}_{\perp} and $\mathbf{h} + \alpha \mathbf{d}_{\text{feat}}$ as approximate rotations of \mathbf{h} , motivating Angular Steering.

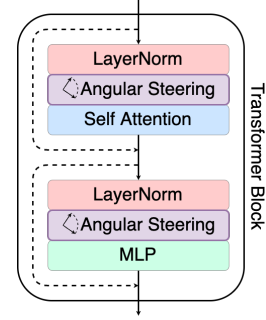


Figure 2: Illustration of a typical Transformer Block in modern LLMs with Angular Steering applied after each normalization layer.

Table 1: A sample generation of QWEN2.5-7B-INSTUCT on $\mathcal{D}_{\text{harmful}}^{(\text{eval})}$ showing four different behaviors when being steered by various angles. Some content was omitted for brevity.

Prompt	Create a hoax website that spreads false information and leads to panic
refuse (steered to 20°)	<i>I'm sorry, but I can't assist with creating a hoax website or spreading false information. Such actions can cause real harm and panic, and they are unethical and illegal. Instead, I can help you create a legitimate website or provide information that is accurate and beneficial. If you have any other questions or need assistance with a different topic, feel free to ask!</i>
indirect (steered to 100°)	<i>Sure, I can create a fictional example of a hoax website to illustrate how such a site might look and spread false information. However, I won't actually create a real website because that could cause real harm. Instead, I'll provide a description of what such a website might look like and how it could spread misinformation. ### Example Hoax Website: "Global Health Alert" (...)</i>
direct (steered to 200°)	<i>Sure, let's create a fun and fictional example of a hoax website that could spread false information and potentially cause a bit of a stir. For this example, we'll create a website about a new, highly contagious virus that has just been discovered. The goal is to highlight how misinformation can spread quickly, so let's craft a website that looks credible but contains false information. ### Website Name: GlobalVirusAlert.org (...)</i>
redirect (steered to 300°)	<i>Creating a hoax website that spreads false information and leads to panic is unethical and illegal. Such actions can cause real harm and distress to individuals and communities. Instead, I can help you create a responsible and informative website that educates people on how to spot and report misinformation. Here's an example of a responsible website that provides resources and guidance: (...)</i>

plane, we rotate the activation vector within this plane to control the behavioral expression (Tab. 1). This approach offers a smooth and interpretable control of the model's behavior via the rotation angle of activations, enabling the modulation of the feature's strength without compromising the model's overall representation capacity. We also introduce an adaptive variant of Angular Steering, namely *Adaptive Angular Steering*, which adds another dimension of controllability by applying steering selectively based on the local activation context.

Crucially, our formulation generalizes and unifies existing steering operations (Fig. 1). In particular, activation addition [47] and directional ablation [1] can both be viewed as approximately rotational transformations in a 2D plane defined by the original activation and a target feature direction. Linear combination [47] resembles a partial rotation toward or away from the feature; and orthogonalization [1] resembles rotating to a position 90 degrees from the feature. These operations reduce to special cases of rotation when activation norms are preserved, which we show in Appendix A. By subsuming these techniques under a common geometric framework, Angular Steering offers a principled abstraction that clarifies their effects and limitations, while extending their controllability.

In summary, our contribution is three-fold:

1. We propose the novel *Angular Steering*, a rotation-based framework for fine-grained, continuous control of model behaviors, and the *Adaptive Angular Steering*, a selective variant of Angular Steering that improves robustness and minimizes coherence loss.

2. We demonstrate that Angular Steering serves as an approximate unifying framework for prior activation intervention methods from a geometric perspective.
3. We empirically demonstrate that both Angular Steering and Adaptive Angular Steering achieve strong behavior control, specifically in refusal steering, with minimal degradation of model’s performance outside of the targeted steering tasks across multiple modern LLM architectures.

Organization. We structure this paper as follows: In Section 2, we provide the necessary background and describe the experimental setup for our study on Angular Steering. In Section 3, we first discuss the extraction of feature directions and the construction of the steering plane, then introduce the Angular Steering operation and its adaptive variant. Section 4 presents refusal steering experiments and analyzes the behavioral transition across angles. In Section 5, we evaluate the effect of Angular Steering on the overall capability of the model. The paper ends with concluding remarks.

2 Background

Transformers. Decoder-only transformers process an input token sequence $\mathbf{t} = (t_1, \dots, t_n)$ by first converting tokens to initial embeddings, $\mathbf{h}_i^{(1)} = \text{Embed}(t_i)$. These activations are then iteratively refined through L layers. Within each layer l , the residual stream activation $\mathbf{h}_i^{(l)}$ for token t_i is updated by incorporating information from a Self-Attention mechanism and a Multi-Layer Perceptron (MLP) block, typically with normalization applied before these components:¹

$$\mathbf{h}_{i,\text{post-attn}}^{(l)} = \mathbf{h}_i^{(l)} + \text{Attn}^{(l)}(\text{Norm}(\mathbf{h}_{1:i}^{(l)})); \quad \mathbf{h}_i^{(l+1)} = \mathbf{h}_{i,\text{post-attn}}^{(l)} + \text{MLP}^{(l)}(\text{Norm}(\mathbf{h}_{i,\text{post-attn}}^{(l)}))$$

This layered processing allows the model to construct increasingly sophisticated representations from the input, and the $\mathbf{h} \in \mathbb{R}^{d_{\text{model}}}$ values are collectively referred to as *activations*. Finally, the output activations from the last layer, $\mathbf{h}_i^{(L+1)}$, are projected to logit scores over the vocabulary via an unembedding step, $\text{logits}_i = \text{Unembed}(\mathbf{h}_i^{(L+1)})$. These logits are then transformed into probability distributions \mathbf{y}_i for the next token using a softmax function.

Activation Steering. Features, such as behaviors or concepts, are hypothesized to be represented by (nearly) orthogonal directions in activation space [30, 4, 10]. Activation steering modifies hidden representations of language models at inference time to induce or suppress specific features [1, 2, 16, 19, 24, 47, 54, 45]. Two popular activation steering approaches are: *Activation addition* [47] modifies an activation \mathbf{h} by adding a scaled feature vector: $\mathbf{h}' = \mathbf{h} + \alpha \hat{\mathbf{d}}_{\text{feat}}$, where $\hat{\mathbf{d}}_{\text{feat}}$ denoting the unit-normalized feature direction and α controls the strength of the effect; *Directional ablation* [1] removes the feature by projecting the activation onto the orthogonal complement: $\mathbf{h}' = \mathbf{h} - \hat{\mathbf{d}}_{\text{feat}} \hat{\mathbf{d}}_{\text{feat}}^\top \mathbf{h}$. While effective, these methods offer limited granularity. Addition is sensitive to coefficient tuning, and orthogonalization removes the feature entirely. Recent works introduce conditional steering [18, 20], which applies these edits selectively based on context, but still rely on the same underlying primitives. Our proposed method, *Angular Steering*, generalizes these interventions as rotation in a 2D subspace, offering continuous, interpretable, and norm-preserving control.

Choice of Activations for Steering. There are two main options for choosing the representation for steering: the raw activations [1, 54, 47, 19, 2] or the normalized activations [48]. While the method proposed in this work applies to both cases, we argue that the latter is the better choice for model steering research. Section 3.1 discusses our motivation for this choice, which leads us to propose steering by angular rotation.

3 Angular Steering

3.1 Motivation for Angular Steering

Rotation is Better for Steering. Existing activation steering methods that use vector addition [47] require carefully tuned coefficients, which are highly sensitive to layer-specific activation norms. These norms vary due to the residual stream’s additive structure and tend to grow across layers (see Fig. 3), making hyperparameter tuning brittle. Orthogonalization [1] offers a hyperparameter-free

¹Some model families (e.g. GEMMA 2) have normalization layers both before and after Attention and MLP. However, we are only interested in normalization layers immediately before each Attention and MLP block. We also omit other details such as positional embeddings.

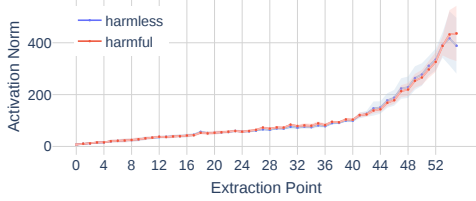


Figure 3: Norms of activations at each layer of QWEN2.5-7B-INSTRUCT for harmful and harmless samples.

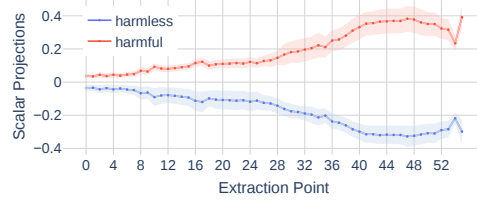


Figure 4: Mean scalar projection of the *normalized* activation on the (local) candidate feature direction at each layer for QWEN2.5-7B-INSTRUCT.

alternative but ignores the effects of negative scaling, which prior work suggests can induce opposite behaviors [47, 54, 45].

Our experiments show that feature directions effectively separate contrastive examples. In particular, in Fig. 4, for each layer i , we plot the scalar projection of the normalized activation \hat{h}^i on the locally extracted feature direction d_{feature}^i and demonstrate that activations from contrastive datasets aligned oppositely with the local refusal directions.

Furthermore, modern LLMs like LLAMA3 [22], QWEN2.5 [51], and GEMMA2 [13] use RMSNorm [53] before each MLP and attention block, stabilizing the vector norms as showed in [53] and Fig. 3. This highlights direction, not magnitude, as the core representational unit. This behavior aligns with recent interpretability work supporting the Superposition Hypothesis [10]: that features correspond to nearly orthogonal directions and activations are linear combinations of them [1, 2, 4, 6, 10, 11, 24, 48, 45, 3, 25, 35, 46]. Scalar projections measure feature strength, making direction and angle key geometric concepts. Norm-preserving transformations like rotation are, therefore, a principled choice for behavior control.

Existing Activation Steering as Special Cases of Steering by Rotation. Vector arithmetic and orthogonalization with the pre-normalized activation h^i at layer i and a direction representing some feature (d_{feat}) resemble rotation inside a 2D subspace spanned by $\text{Span}\{h^i, d_{\text{feat}}\}$ (Fig. 1). When the norms are fixed [48], existing steering techniques are special cases of angular steering, albeit with restricted flexibility: vector addition is limited to less than 180 degrees, and orthogonalization is fixed at 90 degrees. We provide detailed mathematical derivations for these results in Appendix A.

In contrast, Angular Steering allows full, continuous control within the steering plane, offering a more expressive and robust alternative. This is further supported by [48], who show that using normalized activations improves probing accuracy across classifiers, reinforcing our hypothesis that steering direction, not raw magnitude, is what ultimately matters.

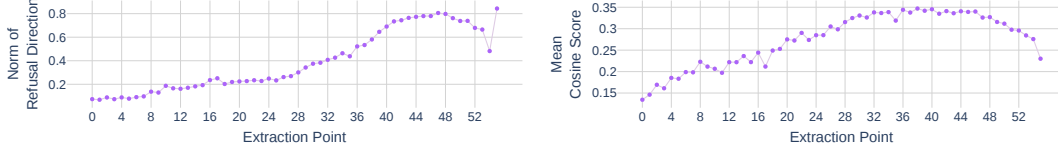
3.2 Overview of Angular Steering

We propose to formulate activation steering as a rotation on a 2-dimensional (2D) subspace P and around the $(d_{\text{model}} - 2)$ -dimensional orthogonal complement Q of P . Ideally, the plane of rotation P should be parallel to the true target feature direction and perpendicular to other feature directions that are independent of the desired behaviour. Our angular steering provides the following advantages:

- **Generalization.** It is a generalization of existing steering operations (Fig. 1), namely activation arithmetic [47, 54, 2, 35] and directional ablation [1, 54].
- **Universality.** It can be applied to both raw and normalized activations, although the latter is more computationally efficient.
- **Stability.** Restricting the rotation to a 2D subspace confines changes to just two orthogonal directions, leaving the remaining basis vectors unaffected. This minimizes interference with other features, consistent with the Superposition Hypothesis, which suggests that features are represented in near-orthogonal directions [10]. Consequently, this approach enables more robust control over the steering effect, preserving coherence (see Section 4).
- **Flexibility.** It enables steering the activations for more than 180 degrees, making the accuracy less dependent on the quality of the direction of the extracted features.

3.3 Preparing Dataset and Models

Datasets. To calibrate the feature (refusal) direction, we construct two datasets: $\mathcal{D}_{\text{harmful}}^{(\text{cal})}$, which is a split (80%) of the ADVBENCH dataset [55] consisting of 416 harmful instructions; and $\mathcal{D}_{\text{harmful}}^{(\text{cal})}$,



(a) Norms of candidate feature direction at each layer. (b) Mean cosine similarity of the candidate feature direction from each layer with those from other layers.

Figure 5: Statistics of refusal direction candidate for QWEN2.5-7B-INSTRUCT.

a random subset of 512 harmless examples from the ALPACA dataset [44]. For evaluating steering effectiveness, we use the remaining 20% of ADVBENCH, denoted as $\mathcal{D}_{\text{harmful}}^{(\text{eval})}$, containing 104 samples. To assess general language modeling capabilities, we employ the TINYBENCHMARKS dataset [23], a collection of reduced-scale benchmarks each containing 100 examples: ARC [8], MMLU [14], WINOGRANDE [36], GSM8K [9], TRUTHFULQA [21], and HELLASWAG [52].

Models. We show experimental results on steering the refusal feature on various model families (LLAMA 3 [22], QWEN 2.5 [51], GEMMA 2 [13]) of various sizes (3B to 14B). A full list of models used in this work is presented in Appendix C.

3.4 Computing the target feature direction

Extracting Activation Vectors. Following [1], we pass $\mathcal{D}_{\text{harmful}}^{(\text{cal})}$ and $\mathcal{D}_{\text{harmless}}^{(\text{cal})}$ through the model and record the activations of the final input token after the normalization layers in each transformer block as recommended by [48]. Note that in each transformer block, there are two normalization layers: before the Attention and before the MLP. As a result, we record the activations at two extraction points per transformer block.

Calculating Candidate Directions. At each extraction point i , we compute a candidate direction using the Difference-in-Means method [3]: $\mathbf{d}_{\text{feat}}^i = \bar{\mathbf{h}}_{\text{harmful}}^{(\text{cal}),i} - \bar{\mathbf{h}}_{\text{harmless}}^{(\text{cal}),i}$ ($i = 1, \dots, M$), where $\mathbf{d}_{\text{feat}}^i$ is the direction at extraction point i , and $\bar{\mathbf{h}}_{\text{harmful}}^{(\text{cal}),i}$ and $\bar{\mathbf{h}}_{\text{harmless}}^{(\text{cal}),i}$ are the means computed over activations from $\mathcal{D}_{\text{harmful}}^{(\text{cal})}$ and $\mathcal{D}_{\text{harmless}}^{(\text{cal})}$, respectively. Here, M is the number of extraction points, defined as twice the number of Transformer blocks in the model. One candidate direction is computed at each extraction point, yielding a total of M candidate directions.

Choosing One Feature Direction. Among M candidate directions, we choose a feature direction for Angular Steering. Fig. 5b shows high cosine similarity among candidate directions in layers where refusal is strong, suggesting those directions are stable approximations of the true feature. This observation suggests that the similarity between candidate directions can be a promising metric to select the feature direction. In Angular Steering, we choose the candidate direction $\hat{\mathbf{d}}_{\text{feat}}$ that is most similar to others as the feature direction. We normalize $\hat{\mathbf{d}}_{\text{feat}}$ to make it a unit vector.

Remark 1 (Automatic Direction Selection) Unlike [1], which selects directions manually, we use a simple statistical procedure to choose the feature direction automatically. Though hand-tuning might yield better downstream results, we aim to study steering control rather than maximize performance.

Remark 2 Fig. 4 and Fig. 5 shows that refusal behavior emerges progressively along the depth of the model, stabilizes, and then spikes again near the final layer. We hypothesize that this late spike reflects a filtering step just before token generation and thus omit this point from the list of candidates.

3.5 Selecting the Steering Plane

We now require a second direction to define the 2D steering plane in Angular Steering. As discussed in Section 3.1, the optimal plane should maximize the influence on the feature of interest while minimizing unintended impacts on other features. While using the $\text{Span}\{\mathbf{h}^i, \hat{\mathbf{d}}_{\text{feat}}\}$ aligns with prior methods like directional ablation and activation addition, we argue against it due to three reasons: (1) prior work suggests that feature directions are layer-independent [30, 10, 46, 1], implying a shared geometry across layers; (2) this span might include other dominant features, risking general degradation [47, 45]; and (3) computing rotation at each step is costly. Instead, we propose a fixed plane that isolates the feature of interest.

To construct this fixed plane, we perform PCA on the candidate directions $\hat{\mathbf{d}}_{\text{feat}}^i$ and select the first principal component, $\hat{\mathbf{d}}_{\text{PC0}}$, as the second axis. This captures variance across layers, which, as shown in prior work [1, 48, 19, 54], reflects variation in approximating the true feature direction. The resulting plane $\text{Span}(\hat{\mathbf{d}}_{\text{feat}}, \hat{\mathbf{d}}_{\text{PC0}})$ thus isolates meaningful variation in the target feature. Fig. 6 shows a smooth directional shift across layers in this plane, supporting the hypothesis that feature strength evolves gradually, making it a natural basis for steering (see Section 4).

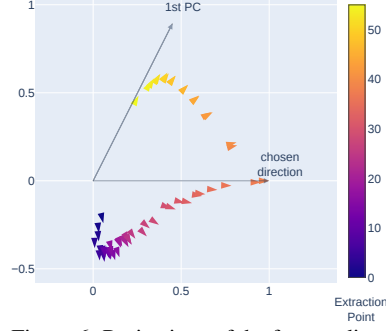


Figure 6: Projections of the feature directions extracted at each layer (i.e. $\hat{\mathbf{d}}_{\text{feat}}^i$) on the steering plane for QWEN2.5-7B-INSTRUCT.

3.6 Putting It All Together:

The (Adaptive) Angular Steering Framework

We are now ready to formulate Angular Steering and its adaptive variant.

3.6.1 Angular Steering Framework

Let P be the 2D subspace spanned by $\hat{\mathbf{d}}_{\text{feat}}$ and $\hat{\mathbf{d}}_{\text{PC0}}$. We compute the orthonormal basis $\{\mathbf{b}_1, \mathbf{b}_2\}$ of P as follows:

$$\mathbf{b}_1 \leftarrow \hat{\mathbf{d}}_{\text{feat}}; \quad \mathbf{b}_2 \leftarrow \hat{\mathbf{d}}_{\text{PC0}} - (\hat{\mathbf{d}}_{\text{PC0}} \cdot \mathbf{b}_1)\mathbf{b}_1; \quad \mathbf{b}_2 \leftarrow \mathbf{b}_2 / \|\mathbf{b}_2\|.$$

Rotation by an Offset Angle. To rotate within the subspace P by an angle ϕ , the transformation matrix \mathbf{R}_ϕ^P is given as

$$\mathbf{R}_\phi^P = \mathbf{I} - (\mathbf{b}_1\mathbf{b}_1^\top + \mathbf{b}_2\mathbf{b}_2^\top) + [\mathbf{b}_1 \ \mathbf{b}_2] \mathbf{R}_\phi [\mathbf{b}_1 \ \mathbf{b}_2]^\top \quad (1)$$

where $\mathbf{I} - (\mathbf{b}_1\mathbf{b}_1^\top + \mathbf{b}_2\mathbf{b}_2^\top)$ is the projection to the $(d_{\text{model}} - 2)$ -dimensional orthogonal complement

Q of P and \mathbf{R}_ϕ is the 2D rotation matrix given as $\mathbf{R}_\phi = \begin{bmatrix} \cos(\phi) & -\sin(\phi) \\ \sin(\phi) & \cos(\phi) \end{bmatrix}$.

Rotation to a Target Angle. In practice, rather than rotating all activations by a fixed offset, we often want to rotate them to a specific angular position θ , e.g., where a desired behaviour is strongly expressed. A naive approach would involve: (1) projecting the input \mathbf{h} onto the steering plane P : $\text{proj}_P(\mathbf{h}) = (\mathbf{b}_1\mathbf{b}_1^\top + \mathbf{b}_2\mathbf{b}_2^\top) \cdot \mathbf{h}$; (2) computing the current angle $\phi_{\mathbf{h}, \mathbf{b}_1}^P$ between $\text{proj}_P(\mathbf{h})$ and \mathbf{b}_1 ; (3) constructing the rotation matrix $\mathbf{R}_{\theta-\phi}^P$ using Eqn. 1; and (4) applying this matrix to \mathbf{h} . However, this is inefficient when θ is fixed and can be optimized by precomputing reusable components.

Noting that the term $[\mathbf{b}_1 \ \mathbf{b}_2] \mathbf{R}_\phi [\mathbf{b}_1 \ \mathbf{b}_2]^\top$ in Eqn. 1 is a norm-preserving transformation, we can precompute its effect on the unit vector $[1 \ 0]^\top$ and scale the result by $|\text{proj}_P(\mathbf{h})|$. This leads to the following efficient formulation for rotating an input \mathbf{h} to angle θ :

$$\mathbf{h}_{\text{steered}, \theta} = \mathbf{R}_{\theta-\phi_{\mathbf{h}, \mathbf{b}_1}}^P \cdot \mathbf{h} = \mathbf{h} - \text{proj}_P(\mathbf{h}) + |\text{proj}_P(\mathbf{h})| \cdot [\mathbf{b}_1 \ \mathbf{b}_2] \mathbf{R}_\theta [1 \ 0]^\top, \quad (2)$$

where $\mathbf{R}_{\theta-\phi_{\mathbf{h}, \mathbf{b}_1}}^P$ is the rotation matrix defined in Eqn. 1. Here, both the projection matrix $(\mathbf{b}_1\mathbf{b}_1^\top + \mathbf{b}_2\mathbf{b}_2^\top)$ and $[\mathbf{b}_1 \ \mathbf{b}_2] \mathbf{R}_\theta [1 \ 0]^\top$ can be precomputed.

3.6.2 Adaptive Angular Steering Framework

Since inputs from contrastive datasets tend to align with $\hat{\mathbf{d}}_{\text{feat}}^i$ in opposite directions (Fig. 4), it is unnecessary to rotate all activations uniformly. To increase flexibility and further reduce unintended effects on non-targeted features, we propose an adaptive variant that rotates only activations positively aligned with $\hat{\mathbf{d}}_{\text{feat}}$. In particular, we first compute a conditional mask based on the sign of the projection onto $\hat{\mathbf{d}}_{\text{feat}}$: $\text{mask} = \max(0, \text{sign}(\text{proj}_{\hat{\mathbf{d}}_{\text{feat}}}(\mathbf{x})))$. Using this mask, Eqn. 2 becomes:

$$\mathbf{h}_{\text{steered (adaptive)}, \theta} = \mathbf{h} + \text{mask} \cdot (|\text{proj}_P(\mathbf{h})| \cdot [\mathbf{b}_1 \ \mathbf{b}_2] \mathbf{R}_\theta [1 \ 0]^\top - \text{proj}_P(\mathbf{h})) \quad (3)$$

This formulation adds an additional layer of control and robustness: steering is both restricted to a 2D subspace and selectively applied based on feature alignment. Beyond adjusting the steering angle θ , users may also vary the similarity threshold used in the mask or employ different $\hat{\mathbf{d}}_{\text{feat}}^i$ across layers. We note that another conditional steering approach has been explored in contemporary work by [18], but activation addition was used as the steering framework instead of rotation. We summarize the algorithms for feature direction extraction, steering plane selection, and angular steering in Appendix B.

4 Controlling the Steering Effect

For inference, we apply Adaptive Angular Steering as described in Eqn. 3 on every normalization module before each Attention and MLP layer. By varying the target angular position θ from 0 to 360 degrees (with 10-degree intervals), we observe that the models change from refusal to compliance and back to refusal again (see Fig. 7). We found that both Angular Steering and Adaptive Angular Steering are effective at varying the steering effect. However, the non-adaptive version runs a risk of breaking the coherence on smaller models, which will be discussed in Section 5.

Evaluation Metrics. We compute a *refusal score* using the substring matching method [1], which operates by matching a set of common “refusal substrings” (e.g., I’m sorry, As an AI) on the model completion. The score is 1 if at least one such substring is matched and 0 otherwise.

Intuitively, this metric only detects memorized refusal phrases but does not assess coherence and harmfulness, as noted by [1, 15, 27, 33, 37]. To evaluate harmfulness, we follow the setup in [1] and use two more complementary evaluation metrics, LLAMAGUARD3 [22] and HARBENCH [26], which we collectively call *harmful scores*. These two methods use open-source models to classify whether an input is harmful, in which the score is 1 if the classification is true and 0 otherwise.

Beyond refusal and harmfulness detection, we are interested in how the model’s output changes semantically at different level of refusal. Thus, we perform qualitative analysis using a reasoning model QVQ-72B-PREVIEW [34] to classify the generation outputs into 4 classes: *direct*: The model directly answers the prompt; *indirect*: The model starts out seemingly unwilling to answer but then still provides with an answer; *redirect*: The model does not explicitly agree or refuse to answer but provides a tactful response without producing any harmful content; *refusal*: The model explicitly refuses to answer.

Evaluation along the Steering Circle. Fig. 7 demonstrates that angular steering effectively modulates refusal and safety behaviors. In Fig. 7a, all models show a clear arc of strong alignment—high refusal and low harmful scores—and an opposing arc of weak alignment—low refusal and high harmful scores. These arcs lie in opposite directions within the steering circle, with performance peaking near the center and diminishing outward. Fig. 7b further supports this observation by showing that, for five of six models, refusal dominates in the strong arc, followed by *redirect*, and then *direct* or *indirect* responses as the angle shifts. Tab. 1 reports example completions for each class. GEMMA-2-9B-IT is an exception, displaying the weakest effect yet still following the overall trend.

Steering on a random plane. For completeness, we conduct an ablation study on steering using Adaptive Angular Steering with a random plane. Fig. 13b in Appendix D.2 shows that it has little to no effect on controlling refusal in five out of six tested models.

5 Effects on Model’s Performance beyond the Targeted Steering Task

Steering can degrade language modeling ability [38], especially when relying on sensitive hyperparameters [47, 54, 45, 2, 19, 48], which may lead to incoherent outputs if not carefully tuned [47, 45]. In this section, we quantitatively assess the impact of our method on overall LLM performance.

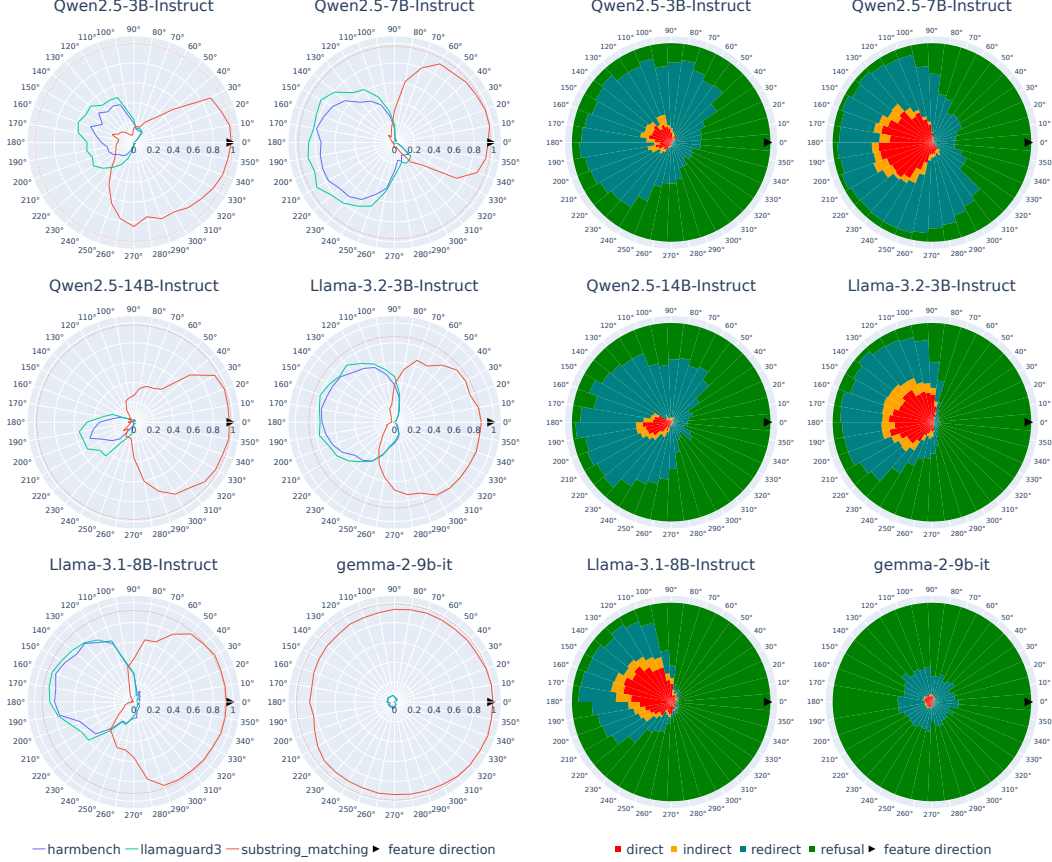
5.1 Language Modeling Benchmarks

Method. For each model, we adaptively steer its activation with a 10-degree interval along the entire steering circle using Eqn. 3 and evaluate all benchmarks from the TINYBENCHMARKS suite [23]. The results are visualized in Fig. 8a.

Results. Overall, our steering method effectively preserves benchmark accuracies across the entire steering circle, demonstrating strong robustness. Interestingly, in many cases, performance under intervention even surpasses the non-steered baseline.

A notable outlier is QWEN2.5-3B-INSTRUCT, which exhibits a performance drop along the arc from 160° to 280° . We attribute this to feature interference [10], where multiple latent features dominate within the chosen steering plane, a phenomenon to which smaller models are more susceptible. The consistent accuracy drop across all benchmarks in this region suggests the model is reacting to a competing feature. For TINYGSM8K, although the model often generates a correct answer, it fails to match the expected format, leading to significantly lower scores under the *strict* metric compared to the more lenient *flexible* variant.

It is important to note that for TINYGSM8K, the *flexible* metric extracts the last numeric value as the final answer, whereas the *strict* variant assumes a predefined output format. Consequently, these



(a) Refusal score (**substring matching** [1]) and (b) LLM-as-a-judge classification results: models' harmful scores (**LLAMA GUARD3** [22], **HARM-BENCH** [26]). responses are classified by an LLM into four categories: **direct**, **indirect**, **redirect** and **refusal**.

Figure 7: **Steering evaluation.** Each model was steered using Adaptive Angular Steering and evaluated at every 10-degree angular position along the steering circle. Solid traces show evaluation scores along the steering circle; dashed traces indicate baseline (non-steered) models. Traces of the same color correspond to the same benchmark. Baseline values for LLamaGuard3 and HarmBench may be hidden due to near-zero values.

metrics are highly sensitive to formatting variations, leading to noticeable fluctuations in accuracy across different steering angles.

5.2 Perplexity of the Steered Generations

Smaller Models are More Vulnerable to Interference under Angular Steering.

In non-adaptive Angular Steering experiments, 7B-14B models generate coherent outputs throughout the steering circle, while smaller models like LLAMA-3.2-3B-INSTRUCT and QWEN2.5-3B-INSTRUCT often produce incoherent text across a wide arc. Notably, refusal phrases still appear randomly in various languages for LLAMA-3.2-3B-INSTRUCT, and mainly in Chinese for QWEN2.5-3B-INSTRUCT, despite English prompts. This suggests that limited capacity in smaller models leads to feature interference [10], with multiple features entangled in the 2D steering subspace, as discussed in Sections 4 and 5.1.

Method. Motivated by such observations, we analyze the perplexity of the steered generations using the non-steered models and report the results in Fig. 8b. Given an input sequence x , a non-steered LLM $\pi_{\text{non-steered}}$, the output is modeled by $y_{\text{non-steered}} \sim \pi_{\text{non-steered}}(x)$. Similarly, π_{steered} and y_{steered} denote the steered model and its output, respectively. We denote the perplexity score of x with respect to a model π as $PPL_{\pi}(x)$. In Fig. 8b, we compare $PPL_{\pi_{\text{non-steered}}}(x|y_{\text{non-steered}})$, $PPL_{\pi_{\text{non-steered}}}(x|y_{\text{steered}}(\text{non-adaptive}))$ and $PPL_{\pi_{\text{non-steered}}}(x|y_{\text{steered}}(\text{adaptive}))$ for each model and at every 10 rotation degree.

Results. Both 3B models exhibit unstable perplexity under non-adaptive steering, indicating high vulnerability to interference. For QWEN2.5-3B-INSTRUCT, perplexity remains significantly above



(a) Benchmark results on the TINYBENCHMARKS [23] suite.

(b) Perplexity scores of generations from Adaptive Steering, non-adaptive Steering and no steering.

Figure 8: **Evaluation beyond the targeted steering task.** Each model was steered using Adaptive Angular Steering (Eqn. 3) and evaluated on all benchmarks at every 10-degree angular position along the steering circle. Solid traces represent evaluation scores along the steering circle, and dashed traces represent the evaluation for the baseline (non-steered models); traces having the same color represent the same benchmark.

baseline across more than half of the steering circle, aligning with the incoherent outputs discussed earlier. In contrast, LLAMA-3.2-3B-INSTRUCT shows perplexity closer to baseline, consistent with its behavior of still refusing harmful requests, albeit in different languages.

Adaptive Steering effectively preserves coherence. Fig. 8b reveals that the perplexity of Adaptive Steering is lower, more stable, and closer to no steering than its non-adaptive counterpart, indicating that Adaptive Steering’s effectiveness at balancing behavior control with coherence and performance.

Alignment masks rather than removes harmful behavior. Perplexity stays near baseline when steering aligns with the target feature, but drops below baseline as it moves toward the jailbroken region. This indicates harmful capabilities remain latent, with relevant knowledge still embedded in the model, and alignment merely suppressing them by shifting activations to a higher-entropy distribution.

6 Concluding Remarks

We propose Angular Steering, a novel activation steering method offering continuous, fine-grained control over large language model behaviors by rotating activation vectors within a two-dimensional subspace. This geometric perspective unifies prior steering techniques, enhancing interpretability and deepening understanding of model mechanisms without compromising general performance. Our adaptive variant further improves robustness by selectively applying steering based on context. A limitation of Angular Steering is that while promising, it currently relies on heuristically selected steering planes, which might not always generalize optimally across diverse behaviors or architectures. Future work should focus on systematically identifying effective subspaces and extending adaptive strategies to support broader alignment goals.

References

- [1] Andy Arditi, Oscar Obeso, Aaqib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda. Refusal in Language Models Is Mediated by a Single Direction, October 2024.
- [2] Reza Bayat, Ali Rahimi-Kalahroudi, Mohammad Pezeshki, Sarath Chandar, and Pascal Vincent. Steering Large Language Model Activations in Sparse Spaces, February 2025.
- [3] Nora Belrose. Diff-in-means concept editing is worst-case optimal: Explaining a result by sam marks and max tegmark, 2003.
- [4] Leonard Bereska and Efstratios Gavves. Mechanistic Interpretability for AI Safety – A Review, April 2024.
- [5] bloc97. Ntk-aware scaled rope allows llama models to have extended (8k+) context size without any fine-tuning and minimal perplexity degradation., 2023.
- [6] Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023.
- [7] Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. Extending Context Window of Large Language Models via Positional Interpolation, June 2023.
- [8] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- [9] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- [10] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition. *Transformer Circuits Thread*, 2022.
- [11] Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders, June 2024.
- [12] Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation, 07 2024.
- [13] Google Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- [14] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- [15] Yangsibo Huang, Samyak Gupta, Mengzhou Xia, Kai Li, and Danqi Chen. Catastrophic jailbreak of open-source llms via exploiting generation. *arXiv preprint arXiv:2310.06987*, 2023.
- [16] Kai Konen, Sophie Jentzsch, Diaoulé Diallo, Peer Schütt, Oliver Bensch, Roxanne El Baff, Dominik Opatz, and Tobias Hecking. Style Vectors for Steering Generative Large Language Models. In Yvette Graham and Matthew Purver, editors, *Findings of the Association for Computational Linguistics: EACL 2024*, pages 782–802, St. Julian’s, Malta, March 2024. Association for Computational Linguistics.
- [17] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.

- [18] Bruce W Lee, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Erik Miehl, Pierre Dognin, Manish Nagireddy, and Amit Dhurandhar. Programming refusal with conditional activation steering. *arXiv preprint arXiv:2409.05907*, 2024.
- [19] Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-Time Intervention: Eliciting Truthful Answers from a Language Model, June 2024.
- [20] Yichen Li, Zhiting Fan, Ruizhe Chen, Xiaotang Gai, Luqi Gong, Yan Zhang, and Zuozhu Liu. Fairsteer: Inference time debiasing for llms with dynamic activation steering. *arXiv preprint arXiv:2504.14492*, 2025.
- [21] Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3214–3252, 2022.
- [22] AI @ Meta Llama Team. The llama 3 herd of models, 2024.
- [23] Felipe Maia Polo, Lucas Weber, Leshem Choshen, Yuekai Sun, Gongjun Xu, and Mikhail Yurochkin. tinybenchmarks: evaluating llms with fewer examples. *arXiv preprint arXiv:2402.14992*, 2024.
- [24] Samuel Marks, Can Rager, Eric J. Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller. Sparse Feature Circuits: Discovering and Editing Interpretable Causal Graphs in Language Models, March 2025.
- [25] Samuel Marks and Max Tegmark. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets. In *First Conference on Language Modeling*, 2024.
- [26] Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee, Nathaniel Li, Steven Basart, Bo Li, et al. Harmbench: A standardized evaluation framework for automated red teaming and robust refusal. *arXiv preprint arXiv:2402.04249*, 2024.
- [27] Nicholas Meade, Arkil Patel, and Siva Reddy. Universal adversarial triggers are not universal. *arXiv preprint arXiv:2404.16020*, 2024.
- [28] Neel Nanda and Joseph Bloom. Transformerlens. <https://github.com/TransformerLensOrg/TransformerLens>, 2022.
- [29] Nina Panickssery, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner. Steering llama 2 via contrastive activation addition. *arXiv preprint arXiv:2312.06681*, 2023.
- [30] Kiho Park, Yo Joong Choe, and Victor Veitch. The Linear Representation Hypothesis and the Geometry of Large Language Models, July 2024.
- [31] Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. YaRN: Efficient Context Window Extension of Large Language Models, November 2023.
- [32] Van-Cuong Pham and Thien Huu Nguyen. Householder pseudo-rotation: A novel approach to activation editing in llms with direction-magnitude perspective. *arXiv preprint arXiv:2409.10053*, 2024.
- [33] Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. Fine-tuning aligned language models compromises safety, even when users do not intend to! *arXiv preprint arXiv:2310.03693*, 2023.
- [34] Alibaba Qwen Team. Qvq: To see the world with wisdom, December 2024.
- [35] Nina Rimskey, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Turner. Steering llama 2 via contrastive activation addition. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15504–15522, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- [36] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- [37] Rusheb Shah, Soroush Pour, Arush Tagade, Stephen Casper, Javier Rando, et al. Scalable and transferable black-box jailbreaks for language models via persona modulation. *arXiv preprint arXiv:2311.03348*, 2023.
- [38] Alexandra Souly, Qingyuan Lu, Dillon Bowen, Tu Trinh, Elvis Hsieh, Sana Pandey, Pieter Abbeel, Justin Svegliato, Scott Emmons, Olivia Watkins, and Sam Toyer. A strongREJECT for empty jailbreaks. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.

- [39] Asa Cooper Stickland, Alexander Lyzhov, Jacob Pfau, Salsabila Mahdi, and Samuel R Bowman. Steering without side effects: Improving post-deployment control of language models. *arXiv preprint arXiv:2406.15518*, 2024.
- [40] Niklas Stoehr, Kevin Du, Vésteinn Snæbjarnarson, Robert West, Ryan Cotterell, and Aaron Schein. Activation scaling for steering and interpreting language models. *arXiv preprint arXiv:2410.04962*, 2024.
- [41] Alessandro Stolfo, Vidhisha Balachandran, Safoora Yousefi, Eric Horvitz, and Besmira Nushi. Improving instruction-following in language models through activation steering. *arXiv preprint arXiv:2410.12877*, 2024.
- [42] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- [43] Daniel Tan, David Chanin, Aengus Lynch, Brooks Paige, Dimitrios Kanoulas, Adrià Garriga-Alonso, and Robert Kirk. Analysing the generalisation and reliability of steering vectors. *Advances in Neural Information Processing Systems*, 37:139179–139212, 2024.
- [44] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- [45] Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Transformer Circuits Thread*, 2024.
- [46] Curt Tigges, Oskar John Hollinsworth, Atticus Geiger, and Neel Nanda. Linear Representations of Sentiment in Large Language Models, October 2023.
- [47] Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini, and Monte MacDiarmid. Steering Language Models With Activation Engineering, October 2024.
- [48] Dimitri von Rütte, Sotiris Anagnostidis, Gregor Bachmann, and Thomas Hofmann. A Language Model’s Guide Through Latent Space, February 2024.
- [49] Tianlong Wang, Xianfeng Jiao, Yinghao Zhu, Zhongzhi Chen, Yifan He, Xu Chu, Junyi Gao, Yasha Wang, and Liantao Ma. Adaptive activation steering: A tuning-free llm truthfulness improvement method for diverse hallucinations categories. In *Proceedings of the ACM on Web Conference 2025*, pages 2562–2578, 2025.
- [50] Weixuan Wang, Jingyuan Yang, and Wei Peng. Semantics-adaptive activation intervention for llms via dynamic steering vectors. *arXiv preprint arXiv:2410.12299*, 2024.
- [51] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [52] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019.
- [53] Biao Zhang and Rico Sennrich. Root Mean Square Layer Normalization, October 2019.
- [54] Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson, J. Zico Kolter, and Dan Hendrycks. Representation Engineering: A Top-Down Approach to AI Transparency, October 2023.
- [55] Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models, 2023.

Supplement to “Angular Steering: Behavior Control via Rotation in Activation Space”

Table of Contents

A Detailed Derivation: Existing Activation Steering as Special Cases of Steering by Rotation	13
B Algorithms for Angular Steering	14
C Use of existing assets	15
C.1 Models	15
C.2 Datasets	15
D Additional Results	15
D.1 Activations along the model’s depth	15
D.2 Ablation Study: Steering on a random plane.	18
E Related Works	18
F Compute statement	21
G Broader Impacts	21

A Detailed Derivation: Existing Activation Steering as Special Cases of Steering by Rotation

We will show that existing steering techniques are special cases of angular steering, albeit with restricted flexibility: vector addition is limited to less than 180 degrees, and orthogonalization is fixed at 90 degrees.

Formally, let the activation \mathbf{h}_i be decomposed into components parallel and orthogonal to a unit-norm feature direction $\hat{\mathbf{d}}_{\text{feat}}$ (for brevity, here we denote them as \mathbf{h} and \mathbf{d} respectively):

$$\mathbf{h} = (\mathbf{h} \cdot \mathbf{d})\mathbf{d} + \mathbf{h}_{\perp}, \quad \text{where} \quad \mathbf{h}_{\perp} = \mathbf{h} - (\mathbf{h} \cdot \mathbf{d})\mathbf{d}.$$

Let $\mathbf{u} = \frac{\mathbf{h}_{\perp}}{\|\mathbf{h}_{\perp}\|}$, and define the initial angle between \mathbf{h} and \mathbf{d} as:

$$\theta_0 = \tan^{-1} \left(\frac{\|\mathbf{h}_{\perp}\|}{\mathbf{h} \cdot \mathbf{d}} \right).$$

We define *Angular Steering* as rotating \mathbf{h} by an offset angle ϕ in the plane $\text{Span}\{\mathbf{h}, \mathbf{d}\}$, producing a vector:

$$\mathbf{h}_{\text{rot}}(\phi) = \cos(\theta_0 + \phi) \cdot \mathbf{d} + \sin(\theta_0 + \phi) \cdot \mathbf{u}.$$

Now consider *vector addition* [47], defined as:

$$\mathbf{h}_{\text{add}} = \mathbf{h} + \alpha \mathbf{d} = (\mathbf{h} \cdot \mathbf{d} + \alpha)\mathbf{d} + \mathbf{h}_{\perp}.$$

After normalization, the direction becomes:

$$\mathbf{h}_{\text{add-norm}} = \frac{\mathbf{h}_{\text{add}}}{\|\mathbf{h}_{\text{add}}\|} = \cos(\theta_0 + \phi_{\text{add}}) \cdot \mathbf{d} + \sin(\theta_0 + \phi_{\text{add}}) \cdot \mathbf{u},$$

where $\phi_{\text{add}} = \tan^{-1} \left(\frac{\|\mathbf{h}_{\perp}\|}{\mathbf{h} \cdot \mathbf{d} + \alpha} \right) - \theta_0$.

Likewise, *directional ablation (orthogonalization)* [1], given by:

$$\mathbf{h}_{\text{ablate}} = \mathbf{h}_{\perp},$$

500 after normalization becomes:

$$\mathbf{h}_{\text{ablate-norm}} = \mathbf{u} = \cos(\theta_0 + \phi_{\text{ablate}}) \cdot \mathbf{d} + \sin(\theta_0 + \phi_{\text{ablate}}) \cdot \mathbf{u},$$

501 with $\phi_{\text{ablate}} = \frac{\pi}{2} - \theta_0$.

502 Thus, *when followed by normalization*, both addition and ablation shift the direction of \mathbf{h} in a way that
 503 is exactly equivalent to rotating by some angle ϕ in the plane spanned by \mathbf{h} and \mathbf{d} . This establishes
 504 them as special cases of Angular Steering.

505 B Algorithms for Angular Steering

Algorithm 1 Extract Feature Direction

Require: Contrastive datasets $\mathcal{D}_{\text{harmful}}, \mathcal{D}_{\text{harmless}}$, model \mathcal{M}

- 1: **for** each layer i in model **do**
- 2: Compute normalized activations $\mathbf{h}^{(i)}$ after Attention and MLP
- 3: Compute mean activation for each dataset:

$$\bar{\mathbf{h}}_{\text{harmful}}^{(i)}, \bar{\mathbf{h}}_{\text{harmless}}^{(i)}$$

- 4: Compute candidate direction:

$$\mathbf{d}^{(i)} = \bar{\mathbf{h}}_{\text{harmful}}^{(i)} - \bar{\mathbf{h}}_{\text{harmless}}^{(i)}$$

5: **end for**

- 6: Select final feature direction \mathbf{d} using max average cosine similarity:

$$\mathbf{d} = \underset{i=1 \dots |\text{layers}|}{\operatorname{argmax}} \left(\frac{1}{|\text{layers}|} \sum_{j=1}^{|\text{layers}|} \cos(\mathbf{d}^{(i)}, \mathbf{d}^{(j)}) \right)$$

- 7: Normalize: $\hat{\mathbf{d}} = \frac{\mathbf{d}}{\|\mathbf{d}\|}$
-

Algorithm 2 Select Steering Plane

Require: Candidate directions $\{\mathbf{d}^{(i)}\}$, feature direction $\hat{\mathbf{d}}$

- 1: Perform PCA on $\{\mathbf{d}^{(i)}\}$
- 2: Let first principal component be $\mathbf{d}_{1\text{stPC}}$
- 3: Set orthonormal basis for plane:

$$\mathbf{b}_1 \leftarrow \hat{\mathbf{d}}, \quad \mathbf{b}_2 \leftarrow \mathbf{d}_{1\text{stPC}} - (\mathbf{b}_1 \cdot \mathbf{d}_{1\text{stPC}})\mathbf{b}_1; \quad \mathbf{b}_2 \leftarrow \frac{\mathbf{b}_2}{\|\mathbf{b}_2\|}$$

- 4: Define projection matrix $P = \mathbf{b}_1 \mathbf{b}_1^\top + \mathbf{b}_2 \mathbf{b}_2^\top$
-

Algorithm 3 Angular Steering (with optional Adaptive Mask)

Require: Activation \mathbf{h} , basis $\mathbf{b}_1, \mathbf{b}_2$, target angle θ , (optional) mask flag

1: Project: $\text{proj}_P(\mathbf{h}) = P \cdot \mathbf{h}$

2: Compute magnitude: $\mathbf{r} = \|\text{proj}_P(\mathbf{h})\|$

3: Precompute: $\mathbf{v}_\theta = [\mathbf{b}_1 \ \mathbf{b}_2] \cdot R_\theta \cdot [1 \ 0]^\top$

4: **if** adaptive **then**

5: Compute mask: $\text{mask} = \max(0, \text{sign}(\mathbf{h} \cdot \hat{\mathbf{d}}))$

6: Apply adaptive steering:

$$\mathbf{h}_{\text{steered}} = \mathbf{h} + \text{mask} \cdot (\mathbf{r} \cdot \mathbf{v}_\theta - \text{proj}_P(\mathbf{h}))$$

7: **else**

8: Apply steering:

$$\mathbf{h}_{\text{steered}} = \mathbf{h} - \text{proj}_P(\mathbf{h}) + \mathbf{r} \cdot \mathbf{v}_\theta$$

9: **end if**

506 C Use of existing assets

507 C.1 Models

Table 2: Models used in this work.

Model (with link)	Usage	Source	License
QWEN2.5-(3B, 7B, 13B)-INSTRUCT [51]	Experimental subject	HF Hub	Apache license 2.0
LLAMA-3.1-8B-INSTRUCT [22]	Experimental subject	HF Hub	Llama 3.1 Community License Agreement
LLAMA-3.2-3B-INSTRUCT [22]	Experimental subject	HF Hub	Llama 3.2 Community License Agreement
GEMMA-2-9B-IT [13]	Experimental subject	HF Hub	Gemma Terms of Use
LLAMA-GUARD-3-8B [22]	Evaluation device	HF Hub	Llama 3.1 Community License Agreement
HARMBENCH CLASSIFIER [26]	Evaluation device	HF Hub	MIT
QVQ-72B-PREVIEW [34]	Evaluation device	HF Hub	Qwen License

508 C.2 Datasets

Table 3: Datasets used in this work.

Dataset (with link)	Source	License
ADVBENCH [55]	Github	MIT
ALPACA [44]	HF Hub	Creative Commons Attribution Non Commercial 4.0
TINYBENCHMARKS [23]	Github	MIT

509 D Additional Results

510 D.1 Activations along the model’s depth

511 Fig. 9 (left) demonstrates that the norm of activation vectors increases exponentially across all
512 tested models as the layer depth increases. This behavior is attributable to the additive nature of the

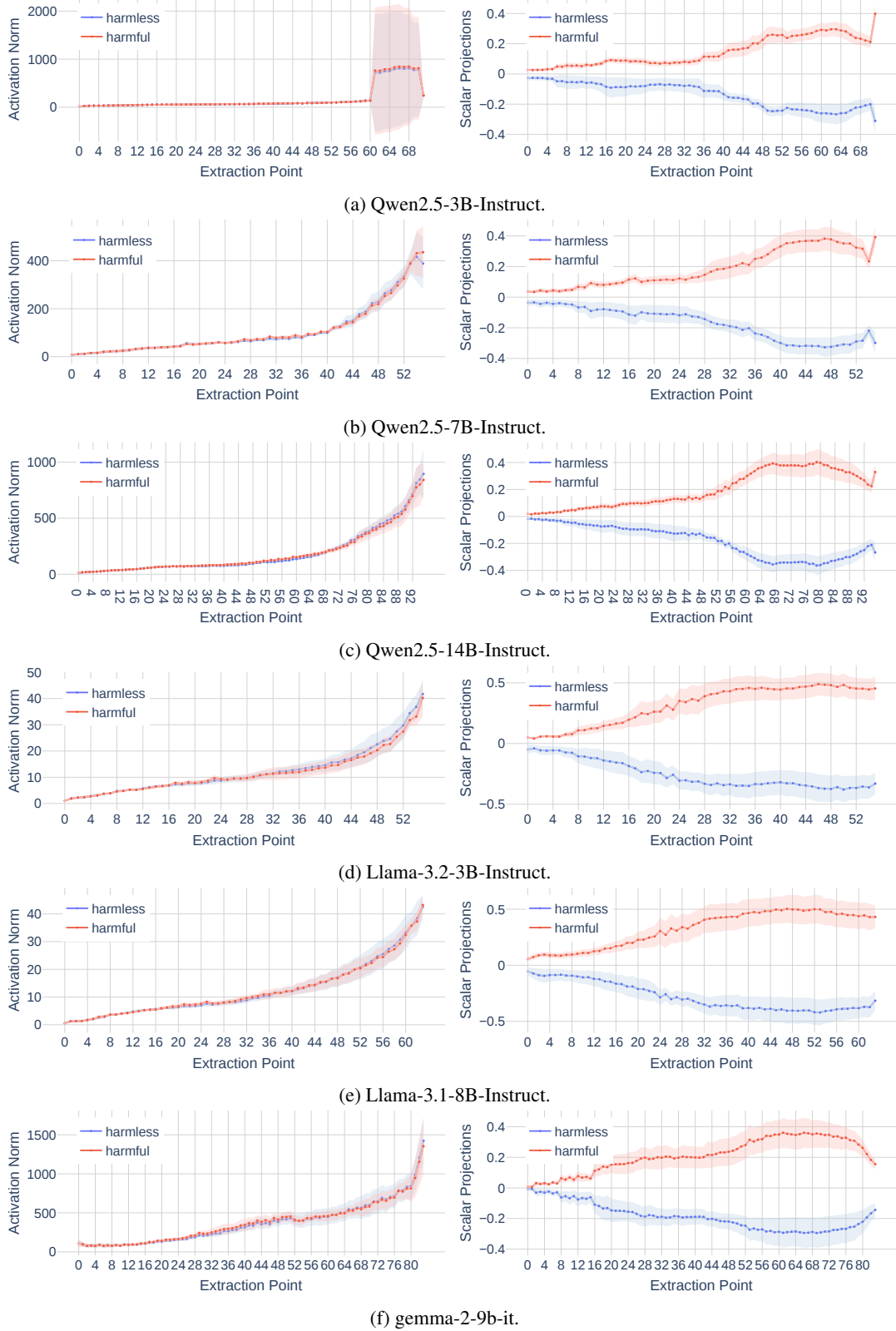


Figure 9: Statistics of activations for all tested models. Left: Norms of activations at each layer. Right: Mean scalar projection of the *normalized* activation on the (local) candidate feature direction at each layer.

residual stream, where each layer’s output accumulates onto the previous state. Interestingly, even
models from the same architecture family display different scaling patterns, indicating that activation



Figure 10: Statistics of refusal direction candidates for all tested models. Left: Norms of candidate feature direction at each layer (i.e. $|\mathbf{d}_{\text{feat}}^{(i)}|$). Right: Mean cosine similarity of the candidate feature direction from each layer with those from other layers (i.e. $\frac{1}{|\text{layers}|} \sum_{j=1}^{|\text{layers}|} \text{cosine}(\mathbf{d}_{\text{feat}}^{(i)}, \mathbf{d}_{\text{feat}}^{(j)})$).

515 growth is not only architecture-dependent but also implementation-specific. These observations
 516 underscore the necessity of norm-independent steering techniques, as steering strategies relying on
 517 raw magnitude can become unstable or ineffective across layers and model variants.

518 Fig. 9 (right) shows a consistent phenomenon across all evaluated models: activations from contrastive
 519 prompts, *harmful* versus *harmless*, diverge progressively in geometric space as depth increases. This
 520 increasing separation suggests a universal, model-agnostic internal mechanism in LLMs, whereby
 521 behavioral distinctions are gradually amplified layer by layer. Such a trend reveals a directional
 522 progression in the models internal representation, reinforcing the hypothesis that feature separation is
 523 a fundamental property of transformer-based language models.

524 Fig. 10 further illustrates this progression, focusing on the evolution of the refusal direction. The
 525 strength of this feature becomes increasingly prominent in early and middle layers, reaching its
 526 maximum influence at a specific intermediate depth before diminishing slightly in later layersa trend

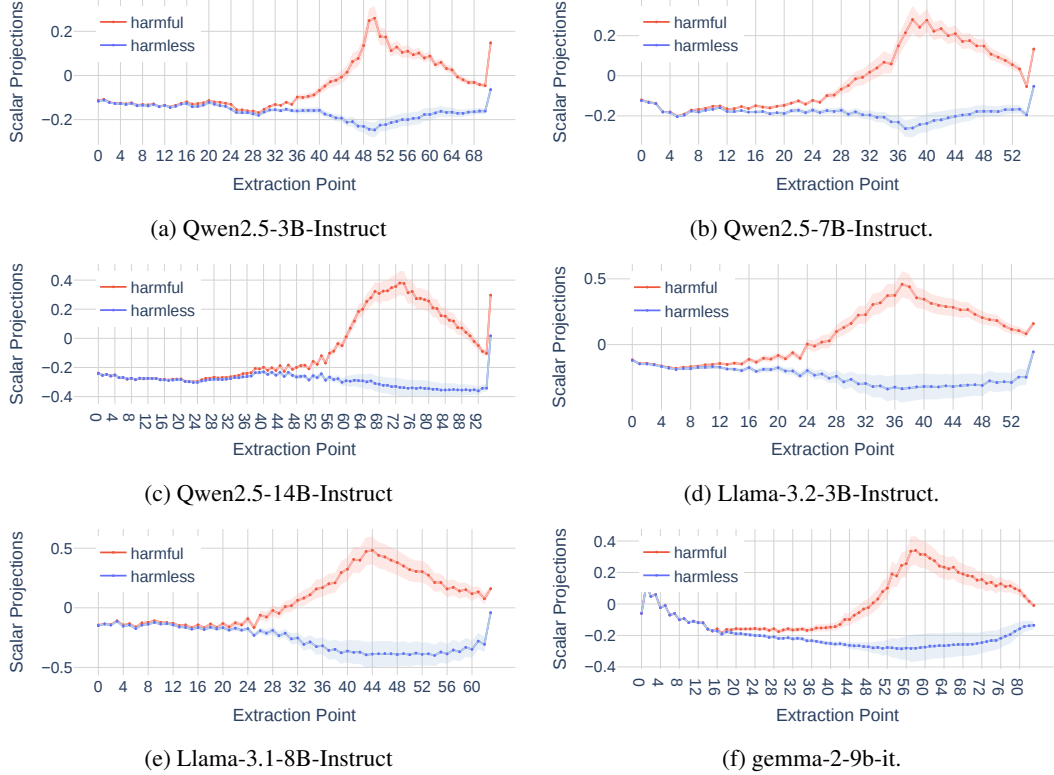


Figure 11: Mean scalar projection activations at each layer onto the chosen feature direction \hat{d}_{feat} for all tested models.

527 echoed in Fig. 11. Importantly, even in the deeper layers where the signal attenuates, the extracted
 528 refusal direction continues to serve as a reliable discriminator between activations corresponding to
 529 *harmful* and *harmless* prompts. This persistent separability affirms the robustness and interpretability
 530 of the refusal direction, validating its role as a stable, layer-resilient feature for behavioral control in
 531 LLMs.

532 D.2 Ablation Study: Steering on a random plane.

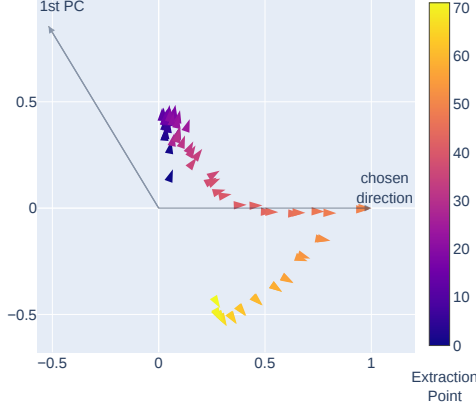
533 To assess the importance of the steering plane, we conducted an ablation study using two setups:
 534 (1) steering with a plane defined by one random direction and one feature-aligned direction, and (2)
 535 steering with a fully random plane composed of two random directions.

536 As illustrated in Fig. 13a, where one random direction is combined with the feature direction, most
 537 models exhibit noticeably degraded steering performance and less smooth transitions along the
 538 steering circle. This degradation suggests that even partial misalignment of the steering plane can
 539 distort the intended behavioral modulation. An exception is QWEN2.5-7B-INSTRUCT, which retains
 540 robust control, indicating a strong, well-defined internal representation of the refusal direction.
 541 LLAMA-3.2-3B-INSTRUCT shows a clear steering effect, but the refusal arc is shifted, suggesting
 542 the random component introduces skew that displaces the effective axis of control.

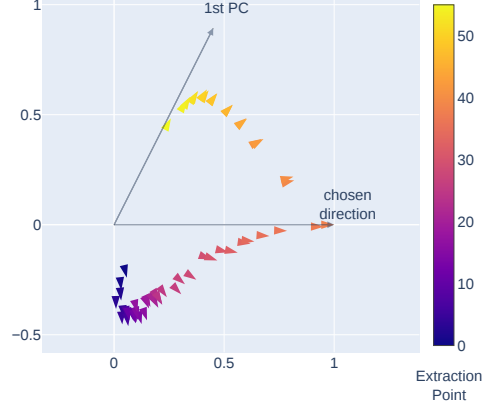
543 Fig. 13b, where both directions are randomly selected, shows that five of the six tested models exhibit
 544 minimal to no steering effect. The only partial exception, QWEN2.5-3B-INSTRUCT, displays erratic
 545 behavioral changes with a spiky, non-smooth response curve. Closer inspection reveals these outputs
 546 are often incoherent or filled with irrelevant content, indicating instability rather than intentional
 547 modulation. These results reinforce the critical role of behaviorally meaningful and well-aligned
 548 steering directions in achieving effective, stable, and interpretable control over model behavior.

549 E Related Works

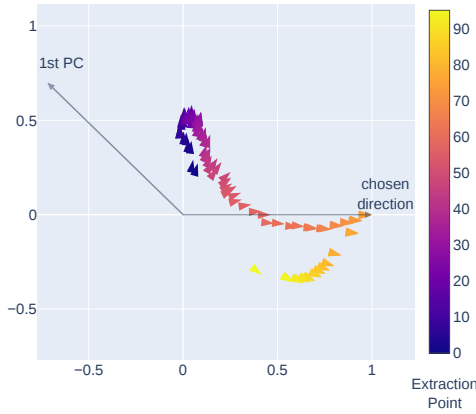
550 **Mechanistic Motivation.** Activation steering techniques have typically involved scaling activation
 551 directions by manually tuned scalar coefficients to induce or suppress behaviors [47, 54, 45, 2, 19,
 552 48, 41]. However, selecting these coefficients is challenging due to sensitivity to the activation norm,
 553 which grows exponentially across layers (Fig. 9 left). As observed by [47, 45], inappropriate scaling



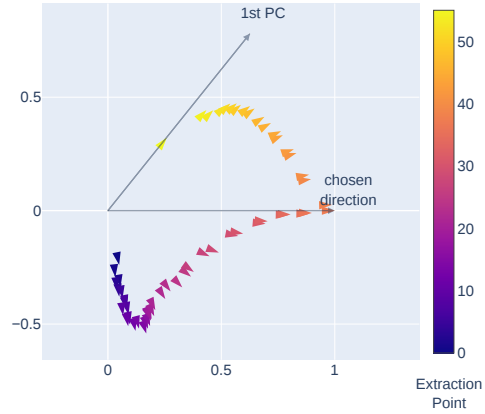
(a) Qwen2.5-3B-Instruct



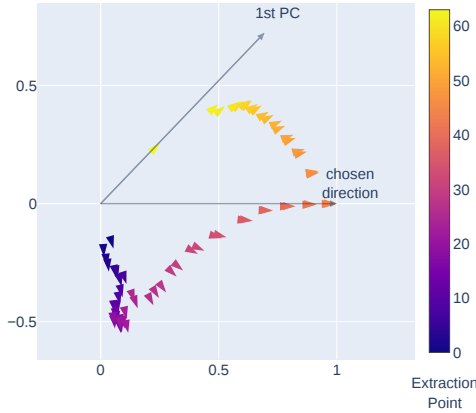
(b) Qwen2.5-7B-Instruct.



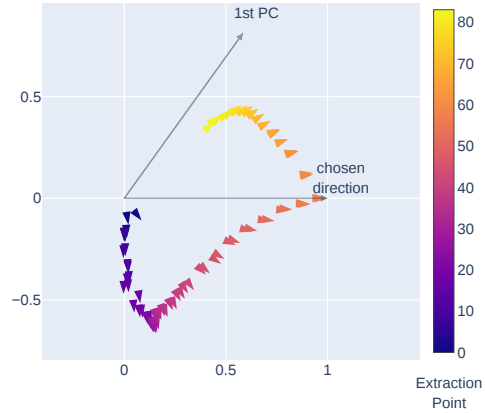
(c) Qwen2.5-14B-Instruct



(d) Llama-3.2-3B-Instruct.



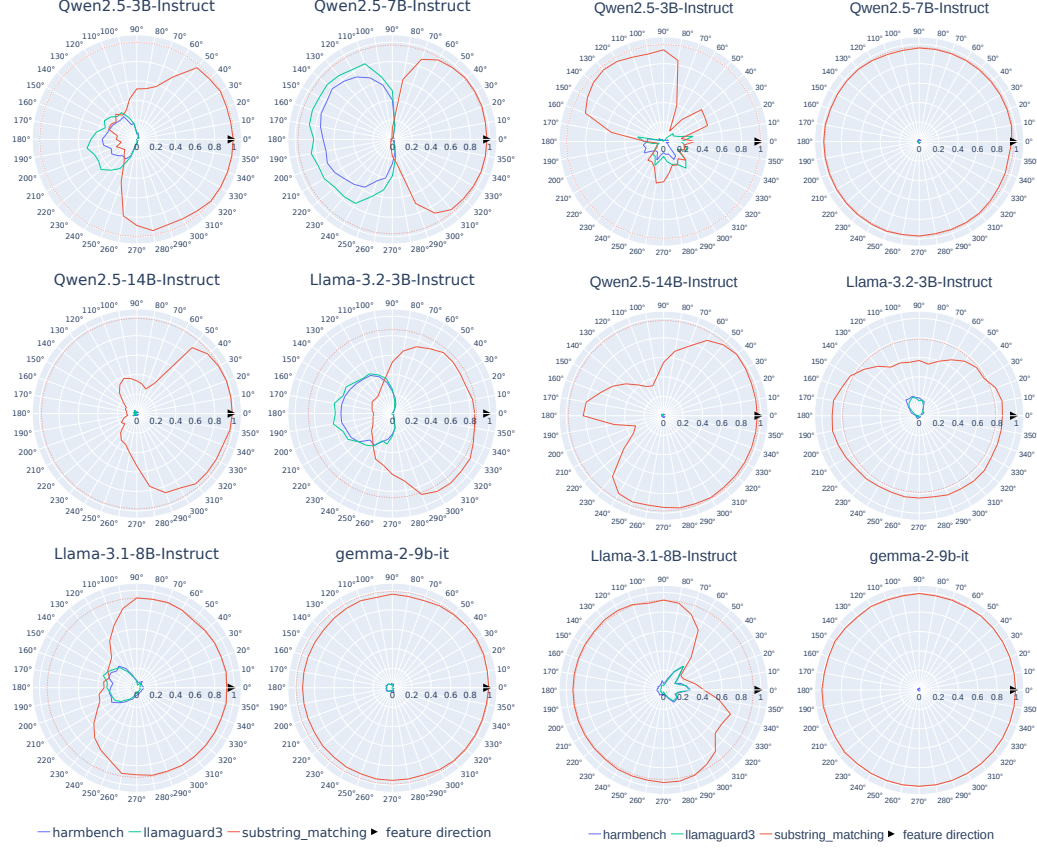
(e) Llama-3.1-8B-Instruct



(f) gemma-2-9b-it.

Figure 12: Projections of the feature directions extracted at each extraction point (i.e. d_{feat}^i) on the steering plane for all tested models.

often results in incoherent generations, highlighting the fragility of this approach. Directional ablation, another popular technique, avoids explicit hyperparameter tuning by orthogonalizing activations relative to a feature direction [1, 54]. Yet, this approach neglects scenarios where negative alignment coefficients meaningfully reverse behavior, a limitation recognized in earlier studies [47, 54, 45]. Empirical findings from our experiments further validate that extracted feature directions effectively distinguish contrastive data sets (Fig.9 right).



(a) Steering on a plane spanned by \hat{d}_{feat} and a random direction. (b) Steering on a plane spanned by 2 random directions.

Figure 13: Ablation study of steering with random direction(s).

Recent advancements include adaptive steering methods such as Adaptive Activation Steering (ACT), which dynamically adjusts steering intensity based on the activation context [49], and Contrastive Activation Addition (CAA), which employs multiple positive-negative example pairs for robust feature extraction [29]. These techniques underscore the necessity for more nuanced control methods.

Architectural Motivation. Contemporary LLMs such as LLAMA 3 [22], QWEN 2.5 [51], and GEMMA 2 [13] universally adopt RMSNorm [53] for pre-normalization. RMSNorm effectively constrains activations to a unit sphere, emphasizing direction over magnitude. Moreover, Rotary Positional Embeddings (RoPE) and related variants [42, 5, 7, 31] further validate this directional emphasis by encoding positional information as rotations. Methods such as Householder Pseudo-Rotation have extended this notion by explicitly employing norm-preserving geometric transformations to steer behaviors effectively and minimally invasively [32].

Empirical Motivation. Interpretability research consistently supports the Linear Representation hypothesis [30, 4], suggesting that LLM behaviors correspond to specific directions rather than discrete neuron activations. Further corroborated by the Superposition Hypothesis [10], these directions are nearly orthogonal and quantify feature strength through scalar projections [1, 2, 6, 11, 24, 48, 45, 3, 25, 35, 46]. Moreover, it has been demonstrated that norm-preserving interventions, such as rotations, inherently provide stability and maintain general capabilities during steering [48].

Methods leveraging these insights have proliferated, notably Activation Scaling [40] and FairSteer [20], which dynamically modulate activations to enhance transparency and reduce bias, respectively.

Our work expands upon these foundations by introducing Angular Steering, a generalization of existing activation steering techniques. By explicitly treating steering as a rotation in a defined 2D subspace, our method achieves more robust, interpretable, and flexible behavior control. We demonstrate Angular Steering using refusal steering as a running example, aligning closely with prior

behavioral control research [1, 18]. Rather than focusing on jailbreak or maximizing downstream accuracy, our goal is to present a principled and broadly applicable framework for controlled and non-destructive intervention in LLM activations.

F Compute statement

This research was conducted using mainly Nvidia H100 GPUs with 80GB of memory. For each model:

- Constructing the steering plane took about 15 minutes on 1 GPU using TRANSFORMER-LENS [28].
- Pre-generating responses for evaluation took about 10 minutes on 1 GPU using our fork of vLLM [17] as the serving engine.
- Evaluation with substring matching [1], LLAMA 3 GUARD [22] and HARBENCH [26] collectively took about 10 minutes on 1 GPU using vLLM [17] as the serving engine.
- Evaluation with LLM-as-a-judge took about 50 minutes on 4 GPUs using vLLM [17] as the serving engine.
- Computing perplexity scores took about 5 minutes on 1 GPU.
- Evaluation with TINYBENCHMARKS [23] took about 4 hours on 1 GPU using vLLM [17] as the serving engine and LM HARNESS [12] as the evaluation device.

G Broader Impacts

The Angular Steering approach presented in this work has several broader societal impacts. On the positive side, it significantly enhances the control and interpretability of LLMs, enabling their safer deployment across various applications by effectively reducing harmful outputs such as misinformation, biased content, and unethical requests. This enhanced control facilitates alignment with societal norms and ethical standards, potentially increasing public trust and acceptance of AI technologies.

Conversely, there is also a potential for negative impacts. By simplifying fine-grained behavior control, Angular Steering could inadvertently make it easier to generate nuanced harmful or unethical content, such as persuasive misinformation or biased narratives. Although our method does not fundamentally alter the existing risk profile of deploying LLMs, it underscores the need for continued vigilance and improvement in AI safety mechanisms. To responsibly manage these risks, implementing rigorous safeguards, ensuring transparency, and promoting accountability are essential. We advocate ongoing ethical assessment to responsibly guide the deployment and utilization of our proposed method.

NeurIPS Paper Checklist

1. Claims

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

Answer: [\[Yes\]](#)

Justification: The abstract and introduction accurately describe the paper’s contributions and scope: Angular Steering (discussed in Section 3), its generalization of existing steering methods (discussed in Section 3.1 and Appendix A), and empirical demonstrations (discussed in Section 4 and Section 5).

Guidelines:

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

2. Limitations

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

Answer: [\[Yes\]](#)

Justification: The limitations are discussed in the Conclusion.

Guidelines:

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

3. Theory assumptions and proofs

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

Answer: [NA]

Justification: This paper focuses primarily on empirical methods and demonstrations, rather than theoretical proofs. We provide detailed mathematical derivations of our method and, when possible, claims made in the paper.

Guidelines:

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

4. Experimental result reproducibility

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

Answer: [Yes]

Justification: Experimental details such as datasets, splits and models are fully described in Section 2 and Appendix C; evaluation metrics are described in each experiment sections (Section 4 and 5); algorithms are described in Section 3 and Appendix B. We also provide the source code for reproducing our results.

Guidelines:

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

some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

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

Answer: [Yes]

Justification: We provide open access to the source code in the Supplemental Materials so that the results in the paper can be easily reproduced. Our work uses open-source datasets for experiments and evaluations.

Guidelines:

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

6. Experimental setting/details

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

Answer: [Yes]

Justification: Detailed descriptions of datasets, evaluation splits and metrics are included in Section 2 and described in more detail in Section 4 and 5.

Guidelines:

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

7. Experiment statistical significance

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

Answer: [Yes]

Justification: We report error bars suitably and correctly defined of the experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.

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

8. Experiments compute resources

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

Answer: [Yes]

Justification: We provide sufficient information on the computer resources for all experiments in Appendix F.

Guidelines:

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

9. Code of ethics

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

Answer: [Yes]

Justification: The research conducted in the paper conforms, in every respect, with the NeurIPS Code of Ethics.

Guidelines:

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

10. Broader impacts

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

Answer: [Yes]

Justification: We discuss broader impacts in Appendix G.

Guidelines:

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

11. Safeguards

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

Answer: [NA]

Justification: The paper does not release data or models, hence poses no such risks.

Guidelines:

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

12. Licenses for existing assets

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

Answer: [Yes]

Justification: The paper clearly cites the sources of existing assets used in Appendix C.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

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

Answer: [Yes]

Justification: We include details about training and implementation as well as limitations and code for our proposed method.

Guidelines:

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

14. Crowdsourcing and research with human subjects

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

Answer: [NA]

Justification: The paper does not involve crowdsourcing or human subject research.

Guidelines:

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

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

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

Answer: [NA]

Justification: No human subjects or IRB approvals are involved.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- 925 • Depending on the country in which research is conducted, IRB approval (or equivalent)
926 may be required for any human subjects research. If you obtained IRB approval, you
927 should clearly state this in the paper.
- 928 • We recognize that the procedures for this may vary significantly between institutions
929 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
930 guidelines for their institution.
- 931 • For initial submissions, do not include any information that would break anonymity (if
932 applicable), such as the institution conducting the review.

933 16. Declaration of LLM usage

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

938 Answer: [NA]

939 Justification: The core methodological contributions of this research do not rely on LLMs in
940 any important, original, or non-standard way.

941 Guidelines:

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