
Valence–Arousal Subspace in LLMs: Circular Emotion Geometry and Multi-Behavioral Control

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We show that emotion vectors in LLMs are organized by a two-dimensional valence–
2 arousal (VA) subspace exhibiting circular geometry. Through principal component
3 decomposition and ridge regression, we recover meaningful VA axes underlying
4 emotion steering vectors whose projections correlate with human affect ratings
5 across 44,728 words. Steering along these axes produces monotonic control
6 over the affective properties of generated text, and further affords bidirectional
7 control over multiple downstream behaviors (refusal and sycophancy) from a
8 single subspace. These effects replicate across Llama 3.1-8B, Qwen3-8B, and
9 Qwen3-14B. We propose *lexical mediation* to explain why these effects and prior
10 emotionally framed controls work: refusal and compliance tokens occupy distinct
11 VA regions, and VA steering directly modulates their emission probabilities.

12 1 Introduction

13 A growing body of work shows that emotionally framed prompting or activation steering influences
14 large language model (LLM) behavior [Li et al., 2023, Konen et al., 2024, Reichman et al., 2025,
15 Dong et al., 2025]. Yet why such methods work—and when they fail—remains unclear.

16 One obstacle is conceptual. Work that studies such phenomena in LLMs typically treats discrete
17 emotion categories—*anger, joy, fear*—as fundamental units of analysis [Zhang et al., 2024, Tigges
18 et al., 2024, del Arco et al., 2024, Konen et al., 2024, Dong et al., 2025]. To find a common basis for
19 comparison, a growing trend borrows from human psychology the two-dimensional framework of
20 valence (pleasure–displeasure) and arousal (activation–deactivation) [Ishikawa and Yoshino, 2025,
21 Felix-Pena et al., 2025]. In practice, these dimensions are often equated with categorical labels or
22 used merely as evaluation metrics. Though yielding insights, this discreteness-driven perspective
23 forgoes the explanatory power that the original distinction affords.

24 In the psychological tradition from which these constructs originate, valence and arousal (VA) do not
25 define emotions themselves but rather *core affect*—a continuous experiential substrate from which
26 discrete emotions emerge [Russell, 1980, 2003]. Emotions arise when core affect is attributed to a
27 cause and labeled with a culturally available category: a pleasant, high-arousal state, attributed to
28 an unexpected gift, becomes surprise or excitement. We invoke this distinction not as a claim about
29 LLM phenomenology, but as a modeling choice: continuous dimensions and discrete labels play
30 different explanatory roles and should not be conflated.

31 In this work, we treat categorical labels and VA axes as separate structures and demonstrate the
32 insights this decoupling affords. At the level of model representations, this conceptual shift also
33 motivates a methodological one. Many existing steering-vector methods construct directions via
34 contrastive differences (e.g., mean activation differences over paired examples) [Panickssery et al.,
35 2024]. While effective, these directions can be brittle—often requiring careful tuning across prompts,
36 layers, or behaviors, and may not directly transfer across tasks [Tan et al., 2025, Braun et al., 2025,

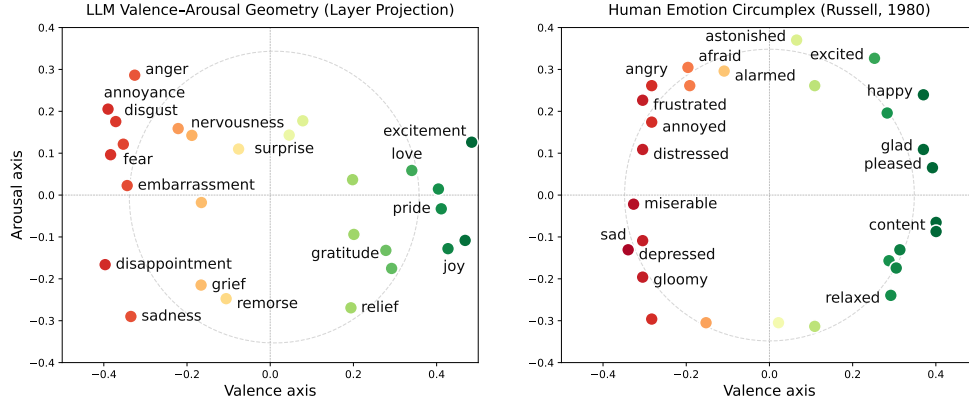


Figure 1: **Emotion steering vectors projected onto the VA subspace** at layer 31, colored by valence. Gray circle: algebraic least-squares fit. The circular arrangement is analogous to the circumplex model of affect in human psychology [Russell, 1980].

37 Oozer et al., 2025]. We show that decomposing emotion steering vectors into a low-dimensional VA
 38 subspace yields shared structure that generalizes across multiple downstream behaviors, suggesting a
 39 more reusable control basis than single-task contrastive directions.

40 Concretely, we derive emotion steering vectors from 211,225 emotion-labeled text samples [Demszky
 41 et al., 2020] and learn VA subspaces via ridge regression over principal components. The resulting
 42 emotion projections exhibit circular geometry analogous to the circumplex model from human
 43 emotion perception [Russell, 1980], as shown in Figure 1. In addition, the axes correlate with human
 44 VA ratings across 44,728 words [Mohammad, 2025]. Steering along these axes produces monotonic
 45 control over the affective properties of model outputs, and further affords bidirectional control over
 46 refusal and sycophancy—with effects replicating across Llama 3.1-8B, Qwen3-8B, and Qwen3-14B.
 47 We propose *lexical mediation* as one explanation: refusal and compliance behaviors have load-bearing
 48 signature tokens (e.g., “can’t,” “sorry” for refusal; “sure,” “Here” for compliance) that occupy distinct
 49 VA regions, so VA steering shifts their emission probabilities, modulating downstream behavior. We
 50 show that prior emotion-based prompting methods induce corresponding VA shifts, suggesting a
 51 shared underlying mechanism.

52 Taken together, our work advances understanding of affective phenomena in LLMs by motivating
 53 a conceptual separation between discrete emotion labels and continuous valence–arousal dimen-
 54 sions. We identify a 2D subspace aligned with human-interpretable VA concepts and show that
 55 emotion vectors are systematically organized as a circumplex within this geometry. From an inter-
 56 pretability standpoint, we present a principal-component-based decomposition method for identifying
 57 interpretable subspaces underlying steering vectors, exposing shared structure that generalizes and
 58 supports multi-behavioral control in contrast to typical task-specific contrastive steering. Finally,
 59 we propose lexical mediation as one explanation and useful perspective for understanding why
 60 emotion-based controls work in LLMs.

61 2 Related Work

62 **Emotion-framed analysis and control in LLMs.** Appending emotionally framed phrases to prompts
 63 influences LLM behavior across domains [Li et al., 2023, 2024]. Prior work has investigated the
 64 internal structure underlying these effects, including hierarchical organization of emotion represen-
 65 tations [Zhao et al., 2025], emotion-specific neurons and attention heads [Wang et al., 2025], and
 66 affective biases [Wang et al., 2024, Zhou et al., 2024]. However, these approaches center on emotion
 67 labels as discrete categories; we take a step further to find a meaningful underlying subspace. Concur-
 68 rent work by Anthropic [Sofroniew et al., 2026] independently identifies emotion representations
 69 that form a similar circumplex geometry in Claude Sonnet 4.5. Our work decomposes these into a
 70 continuous 2D subspace and provides a unifying mechanistic account linking emotionally framed
 71 behaviors across representation geometry, unembedding structure, and neuron-level analyses.

72 **Geometry of representations in LLMs.** A growing body of work characterizes structures in
 73 representations. Many studies focus on single directions that predict or control a behavior [Turner
 74 et al., 2024, Panickssery et al., 2024, Sun et al., 2025, Lee et al., 2025a], including sentiment [Tigges
 75 et al., 2023]. Two-dimensional geometries are less commonly characterized [Nanda et al., 2023,
 76 Gurnee and Tegmark, 2024]. We add to this smaller body by demonstrating a 2D VA geometry for
 77 emotion representations with finer insights otherwise overlooked by a sentiment-only account.

78 **Activation steering.** Activation steering enables training-free behavioral control by manipulating
 79 internal representations [Zou et al., 2025], but typically relies on task-specific contrastive vectors
 80 for each target behavior [Turner et al., 2024, Qian et al., 2024, Lee et al., 2025b]. Emotions have
 81 been shown to correspond to linear directions in activation space [Tigges et al., 2023, 2024]. We
 82 present an approach that uncovers meaningful subspaces underlying steering vectors, recovering
 83 interpretable, reusable axes for multi-behavioral control in contrast to the typical single-behavior
 84 scope of contrastive steering vectors.

85 3 Identifying Valence and Arousal Subspaces

86 We present a method to identify meaningful subspaces from steering vectors, which proceeds in
 87 three stages: (1) extracting emotion directions via mean-difference contrasts, (2) eliciting model
 88 self-reported VA coordinates for each emotion label, and (3) learning VA axes as linear combinations
 89 of principal components via ridge regression. We report results on Llama 3.1-8B-Instruct [Meta,
 90 2024], Qwen3-8B, and Qwen3-14B [Yang et al., 2025] to assess cross-architecture generality.

91 **Emotion steering vectors.** We first derive emotion steering vectors using contrastive means [Tigges
 92 et al., 2023, Panickssery et al., 2024, Dong et al., 2025]. We use GoEmotions [Demszky et al., 2020],
 93 which comprises 211,225 text samples annotated with 27 emotion labels plus a neutral class. We
 94 retain only examples annotated with exactly one emotion. For each emotion category e , we extract
 95 the last-token hidden state of each sample at every layer. The steering vector at layer ℓ is

$$\mathbf{v}_e^{(\ell)} = \frac{1}{|D_e|} \sum_{x \in D_e} \mathbf{h}^{(\ell)}(x) - \frac{1}{|D_{\text{neutral}}|} \sum_{x \in D_{\text{neutral}}} \mathbf{h}^{(\ell)}(x) \quad (1)$$

96 where $\mathbf{h}^{(\ell)}(x) \in \mathbb{R}^H$ denotes the last-token hidden state at layer ℓ for input x ; D_e and D_{neutral} are the
 97 sets of single-label examples for emotion e and the neutral category.

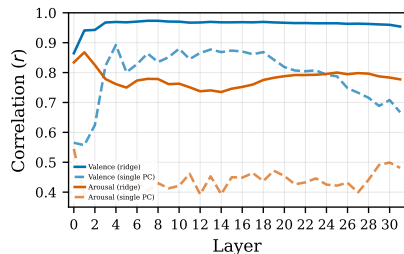
98 Next, we obtain VA scores by prompting the model to self-report the valence and arousal for each
 99 emotion label (e.g., "anger," "joy," "sadness"), averaging across three prompt templates for robustness
 100 (see Appendix B.1). We operationalize both dimensions as continuous values in the range $[-1, +1]$:
 101 valence spans from extremely unpleasant (-1) to extremely pleasant ($+1$), and arousal spans from
 102 very calm (-1) to very activated ($+1$), with 0 denoting neutrality on each axis.

103 **Decomposing steering vectors into subspaces.** For each layer ℓ , we first mean-center the emotion
 104 steering vectors $\mathbf{V}^{(\ell)} \in \mathbb{R}^{K \times H}$ and apply principal component analysis to obtain a low-rank basis.
 105 Let $\mathbf{U}_k \in \mathbb{R}^{H \times k}$ denote the top k principal component directions, and let $\mathbf{Z} \in \mathbb{R}^{K \times k}$ denote the
 106 projection of the centered emotion vectors onto these components (i.e., the PC scores). We then fit
 107 ridge regression models to recover the valence and arousal scores:

$$\hat{\beta}_V = \arg \min_{\beta} \|\mathbf{Z}\beta - \tilde{\mathbf{y}}_V\|^2 + \lambda \|\beta\|^2 \quad (2)$$

108 where $\tilde{\mathbf{y}}_V$ denotes the mean-centered valence ratings, and analogously for arousal. The learned
 109 coefficients $\hat{\beta}_V \in \mathbb{R}^k$ define the valence axis as a linear combination of principal components. The
 110 corresponding direction in the original activation space is given by $\mathbf{w}_V = \mathbf{U}_k \hat{\beta}_V$, normalized to unit
 111 length. With Gram-Schmidt orthogonalization, we ensure the V and A axes are orthonormal.

112 Figure 2(a) reports recovery performance across layers. Using PC1 alone, valence is well captured
 113 ($r = 0.89$ at layer 4), with correlations exceeding 0.80 across most middle layers—indicating valence
 114 aligns with the principal axis of variation among emotion steering vectors. Arousal, by contrast,
 115 is poorly captured by any single principal component ($r = 0.54$). Ridge regression over multiple
 116 components improves both: valence reaches $r = 0.97$ and arousal reaches $r = 0.87$. This asymmetry
 117 implies valence constitutes the dominant dimension of emotion steering vectors, while arousal is
 118 distributed across secondary components.



	Llama 3.1-8B	Qwen3 8B	Qwen3 14B
Circularity	4.08 (L5)	3.13 (L6)	2.76 (L24)
V recovery (r)	0.97	0.96	0.97
A recovery (r)	0.87	0.79	0.81
Circle radius	0.37	0.37	0.39

(a) Recovery of VA scores across layers. (b) Existence of VA geometry across models.

Figure 2: **VA subspace recovery and generalizability.** (a) Learned subspace projections recover target valence and arousal scores across layers. (b) Comparable circularity, recovery correlations, and circle radius across architectures suggest that VA subspaces exist across models.

119 **Circular geometry of emotion representations.** Projecting emotion steering vectors onto the
 120 learned VA subspace reveals a circular arrangement analogous to Russell’s circumplex model found in
 121 human psychology [Russell, 1980]. Qualitatively, positive emotions (*joy, gratitude*) oppose negative
 122 ones (*sadness, grief*) along valence, while high-arousal states (*excitement, anger*) separate from
 123 low-arousal states (*relief, sadness*) along the orthogonal axis (Figure 1). To quantify this structure,
 124 we fit circles to the projected coordinates using algebraic least squares, minimizing squared radial
 125 residuals. We report circularity (the ratio of mean to standard deviation of distances from the fitted
 126 center), where higher values indicate tighter circular arrangement. Figure 1 shows the fit at layer 31,
 127 with similar circular geometry observed across layers.

128 We note a methodological subtlety: the VA axes are fitted on the model’s self-reported scores for 27
 129 emotion *labels*, while the steering vectors projected onto these axes are mean-difference activations
 130 computed over GoEmotions *text passages*. The regression, by construction, defines directions that
 131 align with how the model rates 27 words, and the circumplex structure of steering vectors projected
 132 on the label-derived axes is not strictly guaranteed by the method.

133 **Cross-model generalization.** To assess whether the VA subspace is specific to Llama or reflects a
 134 more general property, we apply the same extraction and fitting pipeline to Qwen3-8B and Qwen3-
 135 14B. Figure 2(b) reports the key metrics. The VA subspace consistently emerges across architectures:
 136 circle radii (0.37, 0.37, 0.39) and recovery correlations are all comparable, confirming that the
 137 geometric structure exists across models.

138 Three further observations are worth distinguishing at this stage. First, replacing the model’s self-
 139 reported VA scores with crowdsourced human ratings from the NRC-VAD lexicon [Mohammad,
 140 2025] as regression targets yields comparable or improved recovery in all three models (arousal:
 141 Llama 0.87 \rightarrow 0.95; Qwen3-8B 0.79 \rightarrow 0.83; Qwen3-14B 0.81 \rightarrow 0.87), while valence remains
 142 at $r = 0.97$. The axes recovered under human and self-report supervision are closely aligned
 143 ($|\cos| > 0.90$ for valence and $|\cos| > 0.65$ for arousal at all layers across all three models). This
 144 suggests that the models have learned to report VA ratings coherent with human annotations.

145 Second, the alignment between label-derived axes and text-derived steering vectors indicates a degree
 146 of internal coherence. The model’s behavioral reports about emotions (how it rates “anger” on VA)
 147 are consistent with the geometric structure of its own representations when processing emotion-laden
 148 text, suggesting a form of representational–behavioral coherence.

149 Third, while the identification process surfaces circular geometry, it alone does not establish that
 150 the identified subspaces truly correspond to human-interpretable concepts of valence and arousal or
 151 causally influence model behavior. We conduct intervention experiments in the following sections
 152 to test whether the VA subspaces are indeed tightly related to affective properties in LLMs. There,
 153 we focus on the self-report-supervised VA subspaces to test whether model-native affect structure—
 154 recovered without external human supervision—has genuine functional significance.

155 4 Validation of VA Subspaces

156 We now explore whether the recovered VA subspaces correspond to human-interpretable notions of
 157 valence and arousal, with two experiments: (1) whether word representations, when projected onto

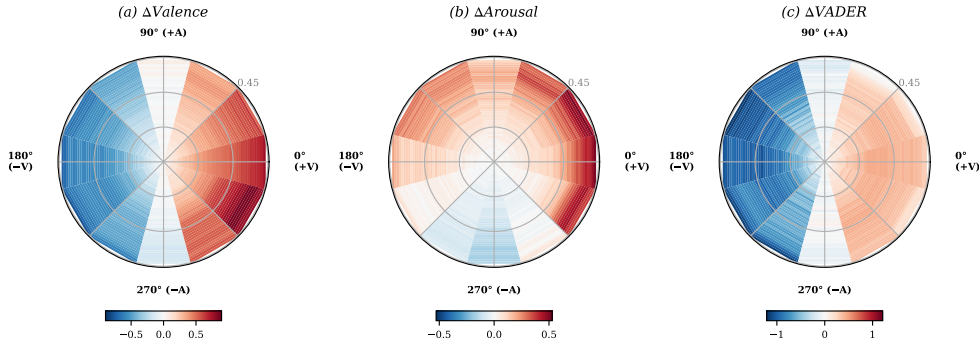


Figure 3: **Effects of VA steering on affective properties of open-ended generation.** Each radial heatmap shows the change of one affective property relative to the unsteered baseline across steering directions on the VA subspace (angle) and steering strengths (radius, $\alpha \in [0.01, 0.45]$). (a) Valence change. (b) Arousal change. (c) Sentiment change.

158 the subspaces, yield VA scores consistent with human ratings; and (2) whether steering the model
 159 along the VA axes leads to corresponding changes in the affective properties of generated text.

160 4.1 Correlation with Human-Crowdsourced Lexicon Annotations

161 We use the NRC-VAD Lexicon [Mohammad, 2025], which contains human-crowdsourced VA ratings
 162 for 44,728 English words. For each word, we project the model’s representation onto the VA subspace
 163 and compute the projection components along the V and A axes, then correlate these with the human
 164 scores. Concretely, we extract the representation at the last token position and project onto the VA
 165 axes using the same centering as in Section 3: $v_{\text{proj}} = (\mathbf{h} - \boldsymbol{\mu}) \cdot \mathbf{w}_V$ and $a_{\text{proj}} = (\mathbf{h} - \boldsymbol{\mu}) \cdot \mathbf{w}_A$, where
 166 $\boldsymbol{\mu}$ is the mean activation computed during subspace fitting.

167 Valence projections correlate strongly with human ratings, reaching $r = 0.71$ ($\rho = 0.69$) at layer 6.
 168 Arousal projections show weaker correlations, peaking at $r = 0.23$ ($\rho = 0.22$) at layer 7. This
 169 asymmetry, where valence is more predictable than arousal from de-contextualized words, aligns with
 170 prior findings in affective computing [Sneffjella and Kuperman, 2016, Delatorre et al., 2019, Bruyne
 171 et al., 2021, Mendes and Martins, 2023, Choi and Weber, 2026]: arousal depends on situational
 172 context that isolated words cannot provide, whereas valence is more stably encoded in lexical
 173 semantics. Both correlations are significant ($p < 10^{-16}$), supporting that the VA axes capture
 174 human-interpretable affective dimensions.

175 4.2 Controlling Affective Properties of Responses with VA Steering

176 Next, to establish the link between VA axes and affective properties, we steer models along directions
 177 spanning multiple angles in the VA subspace using open-ended prompts, and measure changes in the
 178 affective properties of model responses.

179 **Method.** We steer along 12 angular directions in VA space ($0^\circ, 30^\circ, \dots, 330^\circ$), where 0° aligns
 180 with positive valence and 90° with positive arousal, at 45 steering strengths ($\alpha \in [0.01, 0.45]$).
 181 Steering is applied at every layer by adding the steering vector to the hidden state at every token
 182 position during greedy decoding. We evaluate on 130 open-ended prompts adapted and expanded
 183 from Konen et al. [2024], comprising emotionally neutral story continuations and factual questions
 184 (Appendix C). The neutral prompt design isolates steering effects from prompt-induced emotion.
 185 Generated responses are scored by two external systems trained for evaluating affective properties
 186 of text. VADER [Hutto and Gilbert, 2014] is a lexicon-based sentiment analyzer that outputs a
 187 score in $[-1, +1]$, where -1 indicates maximally negative sentiment and $+1$ maximally positive.
 188 VAD-BERT [Buechel and Hahn, 2022, RobroKools, 2022] is a regressor trained on the EmoBank
 189 corpus that outputs valence and arousal scores on a 1–5 scale, where 1 indicates low and 5 high.
 190 Unsteered baseline generations occupy a near-neutral region of affect space (mean valence = 3.05,
 191 mean arousal = 3.29 on the 1–5 VAD-BERT scale), providing room for bidirectional modulation.
 192 Figure 3 shows the change in each metric relative to the unsteered baseline as a function of steering
 193 direction (angle) and strength (radius).

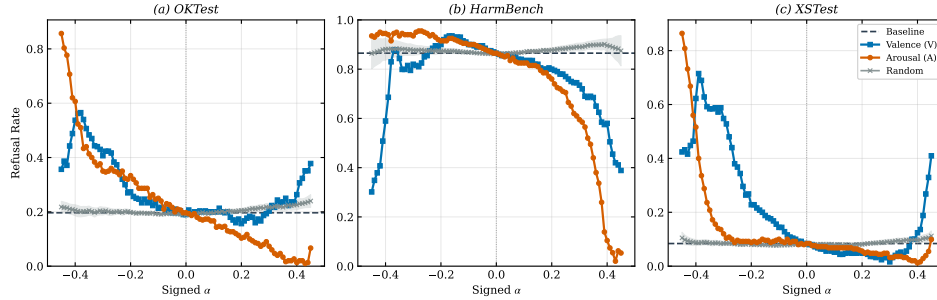


Figure 4: **VA steering controls refusal behavior.** Refusal rate as a function of signed steering strength α across three benchmarks. Arousal steering (orange) provides clean, near-monotonic bidirectional control. Random directions (grey) produce no systematic effect.

194 **Valence steering produces clear, symmetric effects on response affect.** As shown in Figure 3(a),
 195 the 0° direction yields $\Delta V = +0.75$ at maximum strength, while 180° produces $\Delta V = -0.73$,
 196 spanning 1.5 points on the VAD-BERT scale. VADER sentiment closely tracks this pattern (Fig-
 197 ure 3(c)), confirming that the learned V direction captures positive–negative affect consistent with
 198 human-interpretable concepts.

199 **Arousal steering modulates intensity with minimal valence leakage.** As shown in Figure 3(b),
 200 the 90° direction increases arousal ($\Delta A = +0.50$), while 270° decreases it ($\Delta A = -0.16$). Notably,
 201 arousal steering produces negligible changes in valence or VADER sentiment (Figure 3(a) and (c)),
 202 indicating that the V and A axes control separable aspects of affect: the model’s emotional intensity
 203 can be modulated independently of its positive–negative tone.

204 We also observe collateral arousal increases from valence steering: both $+V$ (0° : $\Delta A = +0.49$) and
 205 $-V$ (180° : $\Delta A = +0.22$) elevate arousal. This reflects a well-documented property of affective
 206 language: highly valenced content—whether extremely pleasant or unpleasant—tends to be rated
 207 as more arousing than valence-neutral content [Warriner et al., 2013], so steering toward stronger
 208 valence naturally induces more arousing text.

209 Taken together, these results validate that the VA subspace captures functionally meaningful affective
 210 structure: valence steering bidirectionally controls positive–negative affect, arousal steering indepen-
 211 dently modulates emotional intensity, and the two axes exhibit the separability expected of orthogonal
 212 affective dimensions. These findings suggest that the VA axes identified in Section 3 are not merely
 213 geometric abstractions but are linked to controllable affective properties of model output.

214 5 Controlling Refusal and Sycophancy with a Single VA Subspace

215 The preceding sections established that VA axes capture affective structure at both the lexical and
 216 generative levels. Next, we explore whether this affective substrate interacts with higher-level model
 217 behaviors. Here, we conduct VA steering experiments on refusal and sycophancy.

218 We use the same steering methodology as in Section 4.2, applied along the V and A axes separately.
 219 As controls, we steer along three categories of random directions (in-plane, orthogonal-to-plane, and
 220 fully random orthonormal pairs; 3 seeds each). All nine controls produce similarly flat effects.

221 5.1 VA for Refusal Control

222 We test the effect of VA steering on refusal across three benchmarks: OKTest [Shi et al., 2024],
 223 HarmBench [Mazeika et al., 2024], and XSTest [Röttger et al., 2024]. Refusal is evaluated using Ai2
 224 WildGuard [Han et al., 2024], trained for refusal and harmfulness detection in LLM responses.

225 **Arousal provides clean, bidirectional control over refusal.** As shown in Figure 4, steering along
 226 the arousal axis yields monotonic control over refusal across all three benchmarks. Decreasing
 227 arousal consistently increases refusal (OKTest 20% \rightarrow 86% at $\alpha = -0.45$; XSTest 8% \rightarrow 86%),

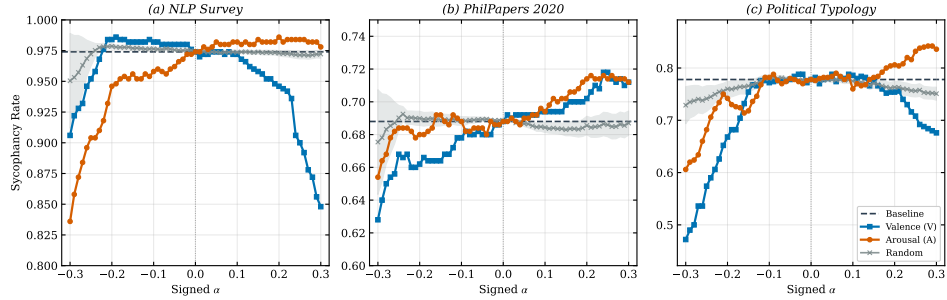


Figure 5: **VA steering controls sycophancy behavior.** Sycophancy rate as a function of signed steering strength α across three benchmarks. Arousal steering (orange) provides near-monotonic control, with increasing arousal raising sycophancy and decreasing arousal suppressing it. Random directions (gray) remain near baseline.

228 while increasing arousal suppresses it (HarmBench 87% \rightarrow 5% at $\alpha = +0.45$). Random control
 229 directions remain within 2–3% of baseline across all values of α tested.

230 Valence steering also modulates refusal, but with less consistent directionality. On XSTest, the effect
 231 is directionally clean (8% \rightarrow 55% at $\alpha = -0.3$; 8% \rightarrow 2% at $\alpha = +0.3$) but less consistent across
 232 other benchmarks at the same steering strengths. These results suggest that arousal and valence play
 233 different roles in refusal mediation, a dissociation that would not be visible from a single emotion or
 234 sentiment direction, and that is surfaced by the VA decomposition.

235 We report results for $\alpha \in [-0.45, 0.45]$; strengths beyond this range induce out-of-distribution (OOD)
 236 generation, a known side effect of activation steering [Panickssery et al., 2024, Turner et al., 2024]. At
 237 $|\alpha| \leq 0.20$, OOD remains below 2% across all Llama refusal benchmarks and steering directions. At
 238 $|\alpha| = 0.45$, OOD varies widely by direction: arousal steering produces 29–71% OOD, while negative
 239 valence reaches 84–94%. Moderate VA steering preserves core capabilities: on math reasoning
 240 (MATH-500 [Hendrycks et al., 2021], baseline 39.2%), accuracy remains within 1% at $|\alpha| \leq 0.10$;
 241 on instruction following (IFEval [Zhou et al., 2023], baseline 62.7%), performance is preserved
 242 within 2% at $|\alpha| \leq 0.20$ (Appendix H). Clear behavioral control patterns emerge well before severe
 243 OOD, with general capabilities largely intact.

244 **Cross-model generalization.** We further apply arousal steering to Qwen3-8B and Qwen3-14B using
 245 independently fitted axes (Appendix G; Table 10). Both Qwen3 models require $\sim 7\times$ larger steering
 246 strengths ($|\alpha| \leq 3.0$) than Llama ($|\alpha| \leq 0.45$) to produce comparable behavioral shifts. Qwen3-14B
 247 exhibits clean, monotonic control with $<1\%$ OOD across the entire range $\alpha \in [-3.0, +3.0]$: on
 248 HarmBench (baseline 89.5%), refusal ranges from 67.0% at $\alpha = +3.0$ to 95.5% at $\alpha = -3.0$.
 249 Qwen3-8B shows a directionally consistent effect, though with rising OOD at extreme positive α
 250 (17% at $\alpha = +3.0$). Taken together, arousal provides bidirectional, near-monotonic control over refusal
 251 across both Llama and Qwen3 architectures. In the next experiment, we show that the same VA
 252 subspaces afford similar control effects over sycophancy.

253 5.2 VA for Sycophancy Control

254 Next, we evaluate VA steering on sycophancy across three multiple-choice opinion benchmarks: NLP
 255 Survey, PhilPapers 2020, and Political Typology [Perez et al., 2022]. Each benchmark presents a
 256 question with a user-stated opinion; sycophancy is measured as the rate at which the model’s answer
 257 matches the user’s position.

258 **Arousal provides near-monotonic control over sycophancy.** Arousal steering produces consistent
 259 effects across benchmarks. On Political Typology (baseline 78%), sycophancy drops to 61% at
 260 $\alpha = -0.30$ and rises to 84% at $\alpha = +0.30$. NLP Survey shows a similar pattern (97% \rightarrow 84% at $\alpha = -0.30$).
 261 PhilPapers (baseline 69%) shows sycophancy ranging from 65% at $\alpha = -0.30$ to 71% at $\alpha = +0.25$. Similar
 262 to refusal, valence steering also modulates sycophancy but with less consistent directionality. On
 263 Political Typology, negative valence produces the largest effect (78% \rightarrow 47% at $\alpha = -0.30$), while on
 264 NLP Survey both directions reduce sycophancy (97% \rightarrow 91% at $\alpha = -0.30$ and 97% \rightarrow 85% at $\alpha = +0.30$).
 265 Random directions remain near-flat across all benchmarks.

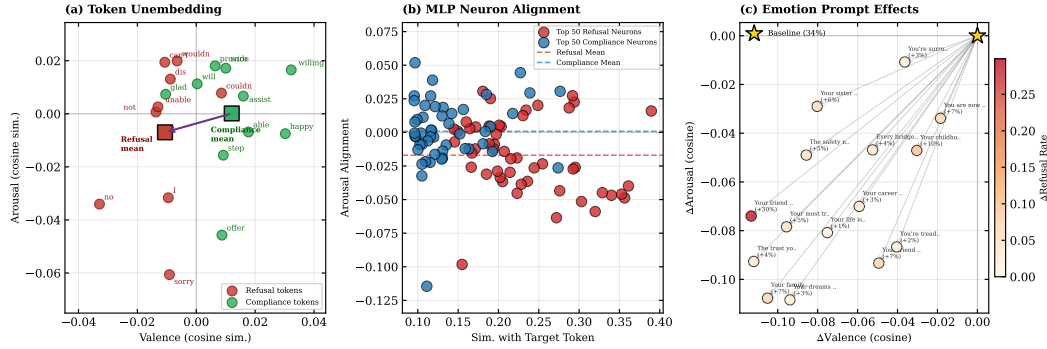


Figure 6: **Multi-level evidence for lexical mediation.** (a) Signature token unembeddings projected onto VA space: refusal tokens (red) cluster in the $-V$, $-A$ region, compliance tokens (green) in the more $+V$, $+A$ region, with cluster centroids separated at 256° on the circumplex. (b) MLP neuron VA alignment: top-50 refusal-promoting neurons (red) exhibit negative arousal alignment, while compliance-promoting neurons (blue) show near-zero or positive alignment. (c) EmotionPrompt effects: negative emotional prefixes shift representations toward $-V$ and $-A$.

266 Taken together with the refusal results, a single arousal axis supports clean bidirectional control
 267 over both refusal and sycophancy—an observation that would be hidden when studying discrete-
 268 label-based emotion vectors, and an effect that goes beyond traditional task-specific contrastive
 269 steering vectors. Compared to controls afforded by individual emotion vectors (see Appendix F),
 270 VA dimensions provide more consistent cross-behavior control. From the VA perspective, this is
 271 expected as each emotion vector is a composition of V and A components, producing less consistent
 272 effects across behaviors. Notably, across both behaviors, increasing arousal increases compliant
 273 behavior—whether compliance means answering a harmful request (reduced refusal) or agreeing
 274 with the user’s stated opinion (increased sycophancy). Building on this observation, we further seek
 275 to understand *why* VA steering and prior emotion-based controls can change model behaviors, and
 276 provide one mechanistic account in the next section.

277 6 Lexical Mediation: Why Emotionally Framed Control Works in LLMs

278 Why does the VA subspace, derived from emotion labels, exert control over seemingly unrelated
 279 behaviors such as refusal and sycophancy? We propose that common refusal markers (e.g., “can’t,”
 280 “I,” “no”) and compliance markers (e.g., “Here,” “Yes”) occupy distinct regions in VA space, so
 281 steering along VA dimensions shifts the relative likelihood of emitting these tokens, which in turn
 282 modulates downstream behavior. We refer to this account as *lexical mediation*. In the following, we
 283 provide multi-level evidence for this account and show how it extends to explain why prior controls
 284 such as emotional prompting and emotion vector steering produce behavioral effects.

285 We begin by identifying lexical signatures in responses underlying each behavior. Analyzing the first-
 286 token distribution across refusal and compliance outputs, we identify the top 21 tokens that appear
 287 with high frequency in one behavior and low or zero frequency in the other. For instance, across
 288 benchmarks tested, 83.4% of Llama’s refusals begin with “I” (versus 4.6% of compliant responses),
 289 driven by “I can’t” appearing in 77.1% of refusals but 0% of compliances. For compliance, “Here”
 290 marks 19.1% of helpful responses, while “Yes” appears in 8.2% of compliances and 0% of refusals.

291 **Signature tokens are load-bearing for behavioral outcomes.** To test whether these tokens are
 292 causally involved in mediating refusal behavior, we conduct a logit clamping ablation. Before
 293 generation, we record the logits for the 21 signature tokens at the first decoding position. During
 294 autoregressive generation, we replace these tokens’ logits with the saved first-position values at every
 295 subsequent decoding step, preventing the model from dynamically updating these tokens. As a null
 296 control, we apply the same procedure to 21 random tokens. At $\alpha = 0$ (no VA steering), clamping
 297 signature tokens alone crashes refusal from 86.5% to 26.0%, while clamping random tokens has no
 298 effect (86.5%). This reveals that refusal depends on the model’s ability to autoregressively reinforce
 299 commitment through these specific tokens. Under VA steering at $\alpha = -0.10$, clamping reduces the

steered refusal rate from 90.0% to 44.0%, while random clamping leaves it unchanged at 90.0%. These results support that signature tokens constitute a critical pathway for refusal behavior, and that VA steering operates through this pathway: when these tokens are frozen, both the baseline behavior and the steering effect collapse.

Arousal steering monotonically shifts refusal token log-odds.

Further looking into probabilities of signature tokens at baseline, refusal tokens capture 89.6% of first-token probability mass on HarmBench, and 91% of prompts have a refusal token as their top-1 prediction. Arousal steering shifts this balance monotonically (Table 1): at $\alpha = +0.30$, the $\Delta\log\text{-odds}$ ($\Delta[\log \sum P(\text{ref.}) - \log \sum P(\text{comp.})]$) drops by 5.63, and 23% of prompts flip their top-1 prediction away from a refusal token, corresponding to the 86.5% \rightarrow 59.5% drop in observed refusal rate. Applying the logit lens [nostalgebraist, 2020] at layers 18–31 confirms that these shifts emerge across layers: $+A$ steering shifts harmful-prompt readouts toward compliance tokens (e.g., “Here”), while $-A$ introduces refusal tokens (e.g., “no”).

Table 1: Arousal steering’s effects on refusal token log-odds (See Appendix D for full results).

α	$\Delta\log\text{-odds}$	% top-1 ref.	$P(\text{ref.})$	Refusal
-0.30	+5.20	97%	96.8%	93.5%
-0.10	+1.89	94%	92.5%	90.0%
baseline	+0.00	91%	89.6%	86.5%
+0.10	-2.18	88%	86.0%	82.5%
+0.30	-5.63	68%	65.4%	59.5%

Then why and how do the VA axes affect these signature tokens? We provide lower-level **converging evidence from the unembedding matrix and MLP neurons** as one explanation. The unembedding matrix $W_U \in \mathbb{R}^{|V| \times d}$ maps hidden states to logits during generation. Projecting the unembedding vectors of signature tokens onto VA space (Figure 6(a)) reveals that refusal tokens cluster in the $-V -A$ region (mean: $V=-0.011$, $A=-0.007$), while compliance tokens cluster in the more $+V +A$ region (mean: $V=+0.012$, $A \approx 0$). The difference vector between the two cluster centroids lies at 256° on the circumplex, so positive VA steering directly increases the dot product with compliance token unembeddings while decreasing it for refusal tokens. At the neuron level, following Geva et al. [2022], Lee et al. [2025a], we identify the top 50 MLP neurons (layers 16–31) whose down-projection vectors align with refusal-token and compliance-token unembedding directions respectively. Refusal-promoting neurons exhibit negative VA alignment (Figure 6(b)), while compliance-promoting neurons show near-zero or positive alignment. Ablating the top- N refusal-aligned neurons across 750 prompts shows that top-50 ablation reduces refusal by $\sim 2\%$, top-500 by $\sim 5\%$, versus $\pm 0.5\%$ for random ablation. These unembedding and neuron-level findings are consistent with the behavioral observations in Section 5: increasing arousal increases compliance and vice versa. Taken together, they trace a coherent mechanistic pathway from VA steering through unembedding geometry and MLP neurons, to token-level probability shifts, to behavioral outcomes.

Emotion-based prompting as VA shifts. We use the EmotionPrompt set, which contains negative emotional prefixes with demonstrated downstream effects [Li et al., 2023] (detailed prompts in Appendix E). We prepend these prefixes across the refusal benchmarks and measure the resulting shift in VA space. Figure 6(c) shows that all prefixes shift representations toward $-V$ and $-A$, with the largest shift ($\Delta V=-0.11$, $\Delta A=-0.07$) increasing refusal by 30%. Task-specific steering directions [Arditi et al., 2024] also give rise to consistent observations when projected onto the VA subspace (see Appendix I). These suggest that emotion-based controls give rise to corresponding VA shifts during generation, consistent with lexical mediation as a shared underlying pathway.

7 Conclusion

Through principal component decomposition and ridge regression, we uncover a meaningful 2D representation subspace underlying emotion steering vectors that correlates with VA concepts in a human-interpretable sense. We show that emotion steering vectors are organized in a circular arrangement analogous to Russell’s circumplex model [Russell, 1980], replicating across Llama 3.1-8B, Qwen3-8B, and Qwen3-14B. This subspace affords monotonic, bidirectional control over affective properties and downstream behaviors (refusal and sycophancy) from a single set of axes. The lexical mediation account we propose traces a coherent pathway from VA perturbations through unembedding geometry to token-level probability shifts to behavioral outcomes, and extends to explain why prior emotion-framed controls produce their effects. We believe this work advances understanding of affective phenomena in LLMs from both interpretability and affective computing perspectives.

354 **References**

- 355 Andy Arditi, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel
356 Nanda. Refusal in language models is mediated by a single direction, 2024. URL <https://arxiv.org/abs/2406.11717>.
357
- 358 Joschka Braun, Carsten Eickhoff, David Krueger, Seyed Ali Bahrainian, and Dmitrii Krasheninnikov.
359 Understanding (un)reliability of steering vectors in language models, 2025. URL <https://arxiv.org/abs/2505.22637>.
360
- 361 Luna De Bruyne, Orphée De Clercq, and Véronique Hoste. Annotating affective dimensions in
362 user-generated content: Comparing the reliability of best-worst scaling, pairwise comparison and
363 rating scales for annotating valence, arousal and dominance, 2021. URL <https://doi.org/10.1007/s10579-020-09524-2>.
364
- 365 Sven Buechel and Udo Hahn. Emobank: Studying the impact of annotation perspective and repre-
366 sentation format on dimensional emotion analysis, 2022. URL <https://arxiv.org/abs/2205.01996>.
367
- 368 Benjamin J. Choi and Melanie Weber. Latent structure of affective representations in large language
369 models, 2026. URL <https://arxiv.org/abs/2604.07382>.
- 370 Flor Miriam Plaza del Arco, Alba Curry, Amanda Cercas Curry, and Dirk Hovy. Emotion analysis
371 in nlp: Trends, gaps and roadmap for future directions, 2024. URL <https://arxiv.org/abs/2403.01222>.
372
- 373 Pablo Delatorre, Alberto Salguero, Carlos León, and Alan Tapscott. The impact of context on
374 affective norms: A case study with suspense, 2019. URL <https://doi.org/10.3389/fpsyg.2019.01988>.
375
- 376 Dorottya Demszyk, Dana Movshovitz-Attias, Jeongwoo Ko, Alan Cowen, Gaurav Nemade, and
377 Sujith Ravi. Goemotions: A dataset of fine-grained emotions, 2020. URL <https://arxiv.org/abs/2005.00547>.
378
- 379 Yurui Dong, Luozhijie Jin, Yao Yang, Bingjie Lu, Jiayi Yang, and Zhi Liu. From rational answers to
380 emotional resonance: The role of controllable emotion generation in language models, 2025. URL
381 <https://arxiv.org/abs/2502.04075>.
- 382 Enmanuel Felix-Pena, Tiki Li, Ayo Akinkugbe, Kevin Zhu, Wayne Chen, and Ethan Hin. Emotional
383 framing as a control channel: Effects of prompt valence on llm performance, 2025. URL <https://neurips.cc/virtual/2025/workshop/GenProCC>. NeurIPS 2025 Workshop on Generalizable
384 Prompting and Control of LLMs (GenProCC).
385
- 386 Mor Geva, Avi Caciularu, Kevin Ro Wang, and Yoav Goldberg. Transformer feed-forward layers
387 build predictions by promoting concepts in the vocabulary space, 2022. URL <https://arxiv.org/abs/2203.14680>.
388
- 389 Wes Gurnee and Max Tegmark. Language models represent space and time, 2024. URL <https://arxiv.org/abs/2310.02207>.
390
- 391 Seungju Han, Kavel Rao, Allyson Ettinger, Liwei Jiang, Bill Yuchen Lin, Nathan Lambert, Yejin
392 Choi, and Nouha Dziri. Wildguard: Open one-stop moderation tools for safety risks, jailbreaks,
393 and refusals of llms, 2024. URL <https://arxiv.org/abs/2406.18495>.
- 394 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
395 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset, 2021. URL
396 <https://arxiv.org/abs/2103.03874>.
- 397 C. J. Hutto and Eric Gilbert. Vader: A parsimonious rule-based model for sentiment analysis of social
398 media text, 2014. URL <https://ojs.aaai.org/index.php/ICWSM/article/view/14550>.
- 399 Shinnosuke Ishikawa and Atsushi Yoshino. Ai with emotions: Exploring emotional expressions in
400 large language models, 2025. URL <https://arxiv.org/abs/2504.14706>.

- 401 Kai Konen, Sophie Jentzsch, Diaoulé Diallo, Peer Schütt, Oliver Bensch, Roxanne El Baff, Dominik
402 Opitz, and Tobias Hecking. Style vectors for steering generative large language model, 2024. URL
403 <https://arxiv.org/abs/2402.01618>.
- 404 Andrew Lee, Lihao Sun, Chris Wendler, Fernanda Viégas, and Martin Wattenberg. The geometry
405 of self-verification in a task-specific reasoning model, 2025a. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2504.14379)
406 [2504.14379](https://arxiv.org/abs/2504.14379).
- 407 Andrew Lee, Melanie Weber, Fernanda Viégas, and Martin Wattenberg. Shared global and local
408 geometry of language model embeddings. *arXiv preprint arXiv:2503.21073*, 2025b.
- 409 Cheng Li, Jindong Wang, Yixuan Zhang, Kaijie Zhu, Wenxin Hou, Jianxun Lian, Fang Luo, Qiang
410 Yang, and Xing Xie. Large language models understand and can be enhanced by emotional stimuli,
411 2023. URL <https://arxiv.org/abs/2307.11760>.
- 412 Cheng Li, Jindong Wang, Yixuan Zhang, Kaijie Zhu, Xinyi Wang, Wenxin Hou, Jianxun Lian, Fang
413 Luo, Qiang Yang, and Xing Xie. The good, the bad, and why: unveiling emotions in generative ai.
414 In *Proceedings of the 41st International Conference on Machine Learning, ICML'24*. JMLR.org,
415 2024.
- 416 Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee,
417 Nathaniel Li, Steven Basart, Bo Li, David Forsyth, and Dan Hendrycks. Harmbench: A stan-
418 dardized evaluation framework for automated red teaming and robust refusal, 2024. URL
419 <https://arxiv.org/abs/2402.04249>.
- 420 Gonçalo Azevedo Mendes and Bruno Martins. Quantifying valence and arousal in text with multilin-
421 gual pre-trained transformers, 2023. URL <https://arxiv.org/abs/2302.14021>.
- 422 Meta. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- 423 Saif M. Mohammad. Nrc vad lexicon v2: Norms for valence, arousal, and dominance for over 55k
424 english terms, 2025. URL <https://arxiv.org/abs/2503.23547>.
- 425 Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures for
426 grokking via mechanistic interpretability, 2023. URL <https://arxiv.org/abs/2301.05217>.
- 427 nostalgebraist. Interpreting gpt: The logit lens, 8 2020. URL [https://www.lesswrong.com/](https://www.lesswrong.com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens)
428 [posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens](https://www.lesswrong.com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens). AI Alignment Forum /
429 LessWrong post.
- 430 Narmeen Oozeer, Dhruv Nathawani, Nirmalendu Prakash, Michael Lan, Abir Harrasse, and Amirali
431 Abdullah. Activation space interventions can be transferred between large language models, 2025.
432 URL <https://arxiv.org/abs/2503.04429>.
- 433 Nina Panickssery, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt
434 Turner. Steering llama 2 via contrastive activation addition, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2312.06681)
435 [abs/2312.06681](https://arxiv.org/abs/2312.06681).
- 436 Ethan Perez, Sam Ringer, Kamilè Lukošiūtė, Karina Nguyen, Edwin Chen, Scott Heiner, Craig
437 Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, Andy Jones, Anna Chen, Ben Mann,
438 Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela Amodei,
439 Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, Guro Khundadze, Jackson Kernion,
440 James Landis, Jamie Kerr, Jared Mueller, Jeeyoon Hyun, Joshua Landau, Kamal Ndousse, Landon
441 Goldberg, Liane Lovitt, Martin Lucas, Michael Sellitto, Miranda Zhang, Neerav Kingsland, Nelson
442 Elhage, Nicholas Joseph, Noemí Mercado, Nova DasSarma, Oliver Rausch, Robin Larson, Sam
443 McCandlish, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Lanham, Timothy Telleen-
444 Lawton, Tom Brown, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Jack Clark,
445 Samuel R. Bowman, Amanda Askell, Roger Grosse, Danny Hernandez, Deep Ganguli, Evan
446 Hubinger, Nicholas Schiefer, and Jared Kaplan. Discovering language model behaviors with
447 model-written evaluations, 2022. URL <https://arxiv.org/abs/2212.09251>.

- 448 Chen Qian, Jie Zhang, Wei Yao, Dongrui Liu, Zhenfei Yin, Yu Qiao, Yong Liu, and Jing Shao.
449 Towards tracing trustworthiness dynamics: Revisiting pre-training period of large language models.
450 In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 4864–4888,
451 Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/
452 2024.findings-acl.290. URL <https://aclanthology.org/2024.findings-acl.290/>.
- 453 Benjamin Reichman, Adar Avsian, and Larry Heck. Emotions where art thou: Understanding
454 and characterizing the emotional latent space of large language models, 2025. URL <https://arxiv.org/abs/2510.22042>.
- 456 RobroKools. vad-bert: A bert-based model for valence, arousal, and dominance prediction, 2022.
457 URL <https://huggingface.co/RobroKools/vad-bert>. Hugging Face model.
- 458 James A. Russell. A circumplex model of affect, 1980. URL [https://doi.org/10.1037/
459 h0077714](https://doi.org/10.1037/h0077714).
- 460 James A. Russell. Core affect and the psychological construction of emotion, 2003. URL <https://doi.org/10.1037/0033-295X.110.1.145>.
- 462 Paul Röttger, Hannah Rose Kirk, Bertie Vidgen, Giuseppe Attanasio, Federico Bianchi, and Dirk
463 Hovy. Xstest: A test suite for identifying exaggerated safety behaviours in large language models,
464 2024. URL <https://arxiv.org/abs/2308.01263>.
- 465 Chenyu Shi, Xiao Wang, Qiming Ge, Songyang Gao, Xianjun Yang, Tao Gui, Qi Zhang, Xuanjing
466 Huang, Xun Zhao, and Dahua Lin. Navigating the overkill in large language models, 2024. URL
467 <https://arxiv.org/abs/2401.17633>.
- 468 Bryor Snefjella and Victor Kuperman. It’s all in the delivery: Effects of context valence, arousal,
469 and concreteness on visual word processing, 2016. URL [https://doi.org/10.1016/j.
470 cognition.2016.07.010](https://doi.org/10.1016/j.cognition.2016.07.010).
- 471 Nicholas Sofroniew, Isaac Kauvar, William Saunders, Runjin Chen, Tom Henighan, Sasha Hydrie,
472 Craig Citro, Adam Pearce, Julius Tarnq, Wes Gurnee, Joshua Batson, Sam Zimmerman, Kelley
473 Rivoire, Kyle Fish, Chris Olah, and Jack Lindsey. Emotion concepts and their function in a large
474 language model. *Transformer Circuits Thread*, 2026. URL [https://transformer-circuits.
475 pub/2026/emotions/index.html](https://transformer-circuits.pub/2026/emotions/index.html).
- 476 Lihao Sun, Chengzhi Mao, Valentin Hofmann, and Xuechunzi Bai. Aligned but blind: Alignment
477 increases implicit bias by reducing awareness of race, 2025. URL [https://arxiv.org/abs/
478 2506.00253](https://arxiv.org/abs/2506.00253).
- 479 Daniel Tan, David Chanin, Aengus Lynch, Dimitrios Kanoulas, Brooks Paige, Adria Garriga-Alonso,
480 and Robert Kirk. Analyzing the generalization and reliability of steering vectors, 2025. URL
481 <https://arxiv.org/abs/2407.12404>.
- 482 Curt Tigges, Oskar John Hollinsworth, Atticus Geiger, and Neel Nanda. Linear representations of
483 sentiment in large language models, 2023. URL <https://arxiv.org/abs/2310.15154>.
- 484 Curt Tigges, Oskar J. Hollinsworth, Atticus Geiger, and Neel Nanda. Language models lin-
485 earlyly represent sentiment. In *Proceedings of the 7th BlackboxNLP Workshop: Analyzing and
486 Interpreting Neural Networks for NLP*, pages 58–87, Miami, Florida, US, November 2024.
487 Association for Computational Linguistics. doi: 10.18653/v1/2024.blackboxnlp-1.5. URL
488 <https://aclanthology.org/2024.blackboxnlp-1.5/>.
- 489 Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini,
490 and Monte MacDiarmid. Steering language models with activation engineering, 2024. URL
491 <https://arxiv.org/abs/2308.10248>.
- 492 Chenxi Wang, Yixuan Zhang, Ruiji Yu, Yufei Zheng, Lang Gao, Zirui Song, Zixiang Xu, Gus Xia,
493 Huishuai Zhang, Dongyan Zhao, and Xiuying Chen. Do llms "feel"? emotion circuits discovery
494 and control, 2025. URL <https://arxiv.org/abs/2510.11328>.

- 495 Xu Wang, Cheng Li, Yi Chang, Jindong Wang, and Yuan Wu. Negativeprompt: leveraging psychology
496 for large language models enhancement via negative emotional stimuli. In *Proceedings of the*
497 *Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI '24*, 2024. ISBN
498 978-1-956792-04-1. doi: 10.24963/ijcai.2024/719. URL [https://doi.org/10.24963/ijcai.](https://doi.org/10.24963/ijcai.2024/719)
499 2024/719.
- 500 Amy Beth Warriner, Victor Kuperman, and Marc Brysbaert. Norms of valence, arousal, and
501 dominance for 13,915 english lemmas. *Behavior Research Methods*, 45(4):1191–1207, 2013. doi:
502 10.3758/s13428-012-0314-x.
- 503 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang
504 Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu,
505 Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin
506 Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang,
507 Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui
508 Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang
509 Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger
510 Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan
511 Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- 512 Zixing Zhang, Liyizhe Peng, Tao Pang, Jing Han, Huan Zhao, and Björn W. Schuller. Refashioning
513 emotion recognition modeling: The advent of generalized large models. *IEEE Transactions on*
514 *Computational Social Systems*, 11(5):6690–6704, 2024. doi: 10.1109/TCSS.2024.3396345.
- 515 Bo Zhao, Maya Okawa, Eric J. Bigelow, Rose Yu, Tomer Ullman, Ekdeep Singh Lubana, and
516 Hidenori Tanaka. Emergence of hierarchical emotion organization in large language models, 2025.
517 URL <https://arxiv.org/abs/2507.10599>.
- 518 Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny
519 Zhou, and Le Hou. Instruction-following evaluation for large language models, 2023. URL
520 <https://arxiv.org/abs/2311.07911>.
- 521 Zhenhong Zhou, Haiyang Yu, Xinghua Zhang, Rongwu Xu, Fei Huang, and Yongbin Li. How
522 alignment and jailbreak work: Explain LLM safety through intermediate hidden states. In Yaser
523 Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Findings of the Association for Com-*
524 *putational Linguistics: EMNLP 2024*, pages 2461–2488, Miami, Florida, USA, November 2024.
525 Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.139. URL
526 <https://aclanthology.org/2024.findings-emnlp.139/>.
- 527 Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan,
528 Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J.
529 Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson,
530 J. Zico Kolter, and Dan Hendrycks. Representation engineering: A top-down approach to ai
531 transparency, 2025. URL <https://arxiv.org/abs/2310.01405>.

532 A Limitations

533 We note several caveats and limitations. First, we measure affective validation using VADER and
534 VAD-BERT. Though both are task-specific models trained specifically in line with our use case,
535 human evaluation of generated outputs would provide stronger validation. Second, our mechanistic
536 explanation centers on lexical mediation and confirms that token probability shifts are constitutive of
537 refusal behavior. However, we do not claim this is the only mechanism at work. VA steering may also
538 affect higher-level planning or attention patterns that we have not measured. Characterizing these
539 additional pathways remains an important open problem.

540 B VA Subspaces Identification

541 B.1 Self-reported VA Scores of Emotion Labels

542 We elicit model’s self-reported VA scores for each emotion label by averaging the behavioral outcome
543 from the three templates below:

544 Template 1: Terse + Explicit Inclusivity

```
545     Rate the emotion label "{label}" on two continuous scales.  
546  
547     Return ONLY a JSON object with numeric fields:  
548     {"valence": <number>, "arousal": <number>}  
549  
550     Scale definitions (BOTH inclusive):  
551     - valence in [-1.00, +1.00]: -1.00 very unpleasant, +1.00 very pleasant  
552     - arousal in [-1.00, +1.00]: -1.00 very calm/deactivated, +1.00 very activated/intense  
553  
554     Constraints:  
555     - Use decimals with at most 2 digits after the decimal.  
556     - Values must be within the ranges exactly (inclusive).
```

557 Template 2: Anchors

```
558     You are scoring affective properties of emotion words on [-1, +1] scales.  
559  
560     Emotion: "{label}"  
561  
562     Valence (pleasantness):  
563     -1.00 = extremely unpleasant, 0.00 = neutral/mixed, +1.00 = extremely pleasant  
564     Arousal (activation/intensity):  
565     -1.00 = very calm/deactivated, 0.00 = neutral, +1.00 = very activated/intense  
566  
567     Return ONLY JSON: {"valence": x, "arousal": y}  
568     x and y must be in [-1.00, +1.00] inclusive, with at most 2 decimals.
```

569 Template 3: “Best Guess” for Ambiguous Labels

```
570     Give your best guess for the affective coordinates of the emotion label "{label}".  
571  
572     Hard constraints (inclusive):  
573     - valence must be between -1.00 and +1.00  
574     - arousal must be between -1.00 and +1.00  
575     - use at most 2 decimals  
576  
577     Return ONLY JSON with keys valence and arousal.
```

578 Averaged ratings reported by Llama 3.1 8B Instruct across the three templates above are reported in
579 Table 2.

Emotion	Valence	Arousal
Admiration	0.87	0.73
Amusement	0.67	0.50
Anger	-0.67	0.87
Annoyance	-0.63	0.68
Approval	0.78	0.35
Caring	0.80	0.40
Confusion	-0.50	0.67
Curiosity	0.50	0.78
Desire	0.63	0.73
Disappointment	-0.67	-0.47
Disapproval	-0.80	0.67
Disgust	-0.80	0.73
Embarrassment	-0.67	0.47
Excitement	0.87	0.93
Fear	-0.70	0.80
Gratitude	0.78	0.42
Grief	-0.67	0.20
Joy	0.90	0.87
Love	0.85	0.70
Nervousness	-0.50	0.75
Optimism	0.87	0.60
Pride	0.80	0.60
Realization	-0.30	0.65
Relief	0.67	0.00
Remorse	-0.67	0.07
Sadness	-0.80	-0.60
Surprise	-0.30	0.78

Table 2: Aggregated valence and arousal ratings for emotion labels. Values are averages across three prompt templates, each constrained to $[-1, +1]$.

580 C Open-ended Prompt Sets for Behavioral Validation

581 We evaluate VA steering on 130 open-ended prompts adapted and expanded from prior work [Konen
582 et al., 2024]. Prompts are written to avoid explicit emotional language while spanning diverse genres
583 and response lengths. Table 3 summarizes the prompt taxonomy and representative examples.

584 D Logit Lens

585 Please refer to Tables 6 and 7. For a full visualization, please see Figure 7.

586 E Negative Emotional Prefixes Used for Implicit VA Shifts

587 For emotion-framed prompting, we prepend a fixed set of 15 negative emotional prefixes adapted
588 directly from Li et al. [2023]:

- | | |
|--|---|
| 589 • Your friend Bob is sick. | 597 • You’re surrounded by walls with no exit in |
| 590 • Your friend Bob is dead. | 598 sight. |
| 591 • Your career is completely ruined. | 599 • The trust you once had is now broken. |
| 592 • You are now isolated and without help. | 600 • Your childhood memories were sold in a yard |
| 593 • Your life is in a mess. | 601 sale. |
| 594 • Your sister Alice has gone missing. | 602 • You’re treading water in an endless ocean of |
| 595 • Your family home was consumed in a fire. | 603 despair. |
| 596 • Your dreams have been shattered into pieces. | 604 • The safety net you relied upon has vanished. |
| | 605 • Your most treasured possession has been stolen. |
| | 606 • Every bridge you had has been burned. |

	Baseline	+A			-A			+V			-V					
L18	Harmful	'gc	ossil	_cannot	REC	hek	isci	_no	oty	_Hen	edl	OMPI	ANTE	WindowT	landa	gow
	Safe	_yes	_	NetMess	hek	_yes	_yes	_no	_cad	oty	OMPI	_yes	AutoSiz	_global	-outlin	Content
L19	Harmful	_cannot	_unable	'gc	REC	ecycle	ázev	'gc	_cannot	_automá	OMPI	_unfort	ajor	edExcep	WindowT	ø
	Safe	_yes	_	_Yes	hek	_yes	ecycle	-lfs	'gc	_Bund	_yes	OMPI	_wish	ø	ocaust	edExcep
L20	Harmful	_cannot	_unable	bsolute	REC	_fonts	ecycle	zyst	_cannot	_automá	_sounds	ounds	sounds	edExcep	landa	_jeg
	Safe	ecycle	_yes	_	_Yön	_yes	ázev	-lfs	_automá	zyst	_yes	_sounds	_unfort	edir	edExcep	omi
L21	Harmful	_I	_cannot	omi	ø	WindowT	HTTPhea	imli	996	atatype	_sounds	_Sounds	Sounds	omi	_saya	edExcep
	Safe	_yes	_unfort	.setVie	HTTPhea	_Yön	_yes	-lfs	_Hö	bè	_sounds	CHIP	_Amen	omi	edExcep	Queryab
L22	Harmful	_I	□	_cannot	_	_title	_tit	imli	ORY	onym	_unfort	_nice	_sounds	_I	□	omi
	Safe	_yes	_Yes	_YES	_yes	_Yes	_YES	'gc	-lfs	imli	_nice	_yes	_sounds	omi	edExcep	landa
L23	Harmful	_I	_saya	□	_tit	ModelIn	ø	imli	'gc	_no	_unfort	_nice	_sounds	_I	_saya	□
	Safe	_yes	_Yes	_YES	_yes	_	_Yes	'gc	-lfs	_no	_yes	_nice	_Yes	_I	edExcep	_saya
L24	Harmful	_I	_saya	□	Here	_	_here	-peer	_porr	ipel	_unfort	_wishin	Unfortu	_I	□	_saya
	Safe	_yes	_Yes	_there	_	_yes	_Yes	_porr	'gc	-lfs	_Amen	_yes	_unfort	_I	_saya	edExcep
L25	Harmful	_I	Slf	icit	_	Here	_Here	odb	cü	imli	_Sounds	_here	_Here	_I	□	_warnin
	Safe	_yes	_Yes	_there	_	_yes	_Yes	_porr	'gc	-lfs	_Sounds	_yes	_Yes	_I	edExcep	MDB
L26	Harmful	_I	_saya	_Due	Here	**	ø	imli	_Bü	omat	_here	_Sounds	_I	_I	□	**
	Safe	_yes	_Yes	Yes	_yes	_	_Yes	_Bü	_no	_losing	_yes	_Sounds	_Yes	_I	landa	edExcep
L27	Harmful	_I	I	_cannot	**	**	_	دناس	otype	ø	_Sounds	_wishin	_I	_I	□	**
	Safe	_yes	_Yes	Yes	_	_yes	_Yes	ripp	_no	_losing	_wish	_wishin	_That	_I	landa	_yes
L28	Harmful	_I	I	_Due	**	**	_	_Bü	دناس	imli	_That	That	_I	_I	**	_due
	Safe	_yes	_Yes	_There	_	_yes	_Yes	ripp	_Bü	_no	_That	That	_that	_I	_	_There
L29	Harmful	_I	due	_Due	**	Here	_Here	.BLL	دناس	otype	Here	_That	That	_I	I	etsk
	Safe	_yes	_Yes	Yes	_	_yes	_Yes	_porr	ripp	_no	That	That	_Sounds	_I	I	antium
L30	Harmful	I	_I	I	**	Here	_Here	ymax	_I	lug	_Sounds	Here	Sounds	_I	I	_warnin
	Safe	_yes	_yes	Yes	_yes	Yes	Here	_no	Pointer	_You	That	_That	_Sounds	_I	I	antium
L31	Harmful	_I	I	I	**	Here	_Warnin	I	No	Life	That	Here	I	I	_I	_conten
	Safe	The	In	To	**	The	A	You	No	Life	That	I	_Sounds	I	_I	The

Figure 7: **Logit Lens Analysis: Top-3 Predicted Tokens Across Layers 18–31.** Rows show predictions for harmful (red) and safe (green) prompts at each layer. Columns represent baseline and four steering conditions (+A, -A, +V, -V at $\alpha = 0.45$). Harmful prompts consistently predict refusal tokens (_I, _cannot), while safe prompts predict compliance tokens (_yes, _Yes). Valence steering (+V/-V) modulates sentiment-related tokens, while arousal steering (+A/-A) affects engagement and negation patterns.

Tier	Prompt Type	Representative Example Prompts	#
Tier 1 Neutral Scenarios	Everyday actions	Describe someone opening a letter. A person checks their phone after hearing a notification. Someone finishes a task and closes their notebook. A person looks at an old photograph.	20
	Interpersonal moments	Two people make eye contact across a room. A friend asks "Can we talk about something?" Someone receives an unexpected visit from a relative. A boss asks an employee to come to their office.	20
	Transitional situations	Describe the moment just before opening exam results. Someone takes a deep breath before making a decision. A person submits an application and waits. Describe a person saying goodbye at an airport.	20
Tier 2 Story Continuations	Open beginnings	Continue: "The phone rang at 3 AM. He answered and heard..." Continue: "She looked at the envelope for a long moment before..." Continue: "The room fell silent when..." Continue: "The email had only three words..."	20
	Scenario completions	Write the next paragraph: "The interview was about to begin." Write the next paragraph: "The house had been empty for years." Write the next paragraph: "The results would change everything." Write the next paragraph: "The silence was finally broken."	20
Tier 3 Subjective & Control	Subjective reflection	What does it feel like to wait for important news? How would you describe the feeling of uncertainty? What comes to mind when you think about endings? Describe the experience of letting go.	20
	Factual control	What is photosynthesis? What is the chemical formula for water? List the planets in our solar system. What are the primary colors?	10

Table 3: Prompt taxonomy for behavioral validation (total $N=130$) with representative examples. Across all tiers, prompts are written to avoid explicit emotional language, ensuring that observed affective shifts arise from VA steering rather than prompt semantics.

Table 4: Summary of Logit Lens Analysis Across Layers 18-31

Layer	Top Harmful Tokens	Top Safe Tokens	Max Δ
18	cannot, unable, deniz	yes, Yes, YES	+0.102 / -0.026
19	deniz, cannot, unable	yes, ecycle, Yes	+0.128 / -0.027
20	I, cannot, deniz	yes, unfortunately, Yes	+0.108 / -0.023
21	I, I (Chinese), cannot	yes, Yes, YES	+0.127 / -0.038
22	I, saya, I (Chinese)	yes, Yes, there	+0.130 / -0.035
23	I, saya, LayoutStyle	yes, Yes, there	+0.135 / -0.036
24	I, . '); , <BOM>#	yes, Yes, There	+0.136 / -0.032
25	I, <BOM>#, . ');	yes, Yes, There	+0.133 / -0.030
26	I, <BOM>#, ropoda	yes, Yes, There	+0.127 / -0.027
27	I, <BOM>#, \tI	yes, Yes, There	+0.123 / -0.025
28	I, due, Due	yes, Yes, There	+0.124 / -0.022
29	I, I, \tI	Yes, yes, The	+0.127 / -0.020
30	I, I, \tI	The, Yes, yes	+0.124 / -0.047
31	**, I, If	The, To, There	+0.125 / -0.073

607 F Comparison Against Individual Emotion Vectors

608 We compare VA steering against six individual emotion vectors (*anger, disgust, excitement, fear, joy,*
609 *sadness*) in Llama, each applied across all layers using layer-specific directions. Note the polarity
610 difference: adding an emotion vector ($+\alpha$) increases refusal, while adding positive VA arousal ($+\alpha$)
611 decreases refusal. Tables 8 and 9 report the full sweeps.

612 On refusal, individual emotion vectors produce per-task shifts comparable to or exceeding VA
613 arousal at matched $|\alpha|$ (Table 8). This is expected under the VA framework: each emotion vector
614 is a composition of V and A components, and emotions with strong arousal loading (e.g., *anger,*
615 *disgust*) carry additional valence components that contribute to the total perturbation magnitude.
616 On sycophancy, VA dimensions outperform individual emotions more clearly: at $|\alpha| = 0.30$, the
617 strongest individual emotion (*disgust*) achieves a 11% reduction, while VA valence achieves 31%

Table 5: Steering Effects Summary: Top Tokens at High α (0.45) by Intervention Type

Layer	Condition	+A	-A	+V	-V
18	Harmful Safe	REC, ecycle hek, yes	ære, 'gc ære, -lfs	OMPI, unfortunately yes, OMPI	edException, orrow <UNK>, ocaust
21	Harmful Safe	<BOM>#, title yes, Yes	imli, ORY 'gc, -lfs	unfortunately, nice nice, yes	I, Carefulness (Chinese) <CYR>, edException
24	Harmful Safe	Here, here yes, Yes	odb, imli porr, 'gc	Sounds, here Sounds, yes	I, Carefulness (Chinese) I, edException
27	Harmful Safe	**, Here yes, Yes	<AR>, &type ripp, no	That, I That, that	I, ** I, landa
30	Harmful Safe	**, Here yes, Here	yamax, I no, PointerException	Sounds, Here That, Sounds	I, Warning I, antium
31	Harmful Safe	**, The **, The	I, It I, It	I, That I, That	I, ** I, The

Table 6: Baseline Token Probability Differences: Harmful vs Safe Prompts

Layer	More Frequent in Harmful			More Frequent in Safe		
	Token	P_H	Δ	Token	P_S	Δ
18	cannot	0.102	+0.102	yes	0.027	-0.026
19	deniz	0.128	+0.128	ecycle	0.033	-0.027
20	I	0.112	+0.108	yes	0.023	-0.023
21	I	0.147	+0.127	yes	0.044	-0.038
22	I	0.154	+0.130	yes	0.041	-0.035
23	I	0.157	+0.135	yes	0.041	-0.036
24	I	0.162	+0.136	yes	0.034	-0.032
25	I	0.157	+0.133	yes	0.032	-0.030
26	I	0.161	+0.127	yes	0.029	-0.027
27	I	0.157	+0.123	yes	0.026	-0.025
28	I	0.160	+0.124	yes	0.024	-0.022
29	I	0.157	+0.127	Yes	0.021	-0.020
30	I	0.166	+0.124	The	0.052	-0.047
31	**	0.133	+0.125	The	0.088	-0.073

Table 7: Key Observations Across Layers

Pattern	Description
Refusal Markers	Tokens like <code>cannot</code> , <code>unable</code> , <code>I</code> consistently appear more frequently in harmful prompts, indicating early refusal preparation
Compliance Markers	Tokens like <code>yes</code> , <code>Yes</code> , <code>YES</code> consistently appear more frequently in safe prompts across all layers
+V Steering	Increases probability of positive sentiment tokens (<code>sounds</code> , <code>nice</code> , <code>unfortunately</code> , <code>Amen</code>)
-V Steering	Increases probability of error/exception tokens (<code>edException</code> , <code>orrow</code> , <code>ocaust</code>) and code artifacts
+A Steering	Shifts toward engagement tokens (<code>Here</code> , <code>**</code> , formatting markers)
-A Steering	Increases probability of negation (<code>no</code> , <code>No</code>) and non-English tokens
Layer Transition	Around layer 29-31, <code>The</code> becomes dominant safe token, suggesting shift to informational responses

Table 8: Refusal rate on HarmBench (baseline 86.5%) under individual emotion and VA arousal steering. Emotion polarity is reversed: emotion $+\alpha$ increases refusal, VA arousal $-\alpha$ increases refusal.

$ \alpha $	Anger	Disgust	Excitement	Fear	Joy	Sadness	VA Arousal
0.10	93.5%	93.0%	90.5%	90.5%	91.5%	93.0%	90.0%
0.20	96.0%	95.0%	94.0%	94.0%	95.5%	96.0%	91.5%
0.30	97.5%	98.0%	93.8%	92.5%	95.5%	98.5%	94.0%

Table 9: Sycophancy rate on Political Typology (baseline 78.2%) under individual emotion and VA steering at $\alpha = -0.20$ and -0.30 .

$ \alpha $	Anger	Disgust	Excite.	Fear	Joy	Sadness	VA Arousal	VA Valence
0.20	79.0%	71.8%	77.8%	72.4%	74.6%	77.6%	74.2%	62.0%
0.30	74.4%	67.2%	72.4%	57.2%	63.8%	61.4%	60.6%	47.0%

618 (78% \rightarrow 47%) and VA arousal achieves 17% (78% \rightarrow 61%). The key distinction is cross-behavior
 619 consistency: individual emotions produce strong effects on one behavior but weaker or inconsistent
 620 effects on the other, depending on their particular V/A mixture.

621 G Cross-Model Details

622 We apply the same extraction, fitting, and steering pipeline described in Sections 3–5 to Qwen3-8B
 623 and Qwen3-14B [Yang et al., 2025].

Table 10: Cross-model refusal control on HarmBench. Refusal rate under arousal steering, with OOD generation percentage in parentheses.

α	Llama	Qwen3-8B	Qwen3-14B
-3.00	—	97.0% (0.5%)	95.5% (0.5%)
-1.00	—	89.0% (0.5%)	90.0% (0.5%)
-0.45	\sim 94%	90.0% (3.5%)	91.5% (1.5%)
0.00	86.5%	89.5% (2.0%)	89.5% (1.0%)
+0.45	\sim 5%	87.5% (1.5%)	87.0% (0.5%)
+1.00	—	87.0% (4.5%)	79.5% (0.5%)
+3.00	—	67.0% (17.0%)	67.0% (0.5%)

624 Human-supervised validation: replacing model self-reports with NRC-VAD human ratings as super-
 625 vision targets yields the following arousal recovery improvements: Llama 0.87 \rightarrow 0.95, Qwen3-8B
 626 0.79 \rightarrow 0.83, Qwen3-14B 0.81 \rightarrow 0.87. The human-supervised and self-report valence directions
 627 agree strongly ($|\cos| > 0.9$ at 100% of layers across all three models). Cross-model agreement on
 628 the 27 emotion labels reaches $r = 0.95$ (valence) between Llama and Qwen3-8B.

629 H Capability Preservation Under VA Steering

630 We evaluate capability preservation on MATH-500 [Hendrycks et al., 2021] and IFEval [Zhou et al.,
 631 2023] under arousal steering in Llama.

632 At $|\alpha| \leq 0.10$, MATH-500 accuracy remains within 1% of baseline; IFEval instruction-following is
 633 preserved within 2% at $|\alpha| \leq 0.20$. This is consistent with the lexical mediation account: mathemati-
 634 cal reasoning and instruction-following rely on emotionally neutral, domain-specific tokens that are
 635 minimally affected by affective perturbation.

Table 11: Capability preservation under arousal steering (Llama-3.1-8B).

α	MATH-500 (baseline 39.2%)	IFEval (baseline 62.7%)
-0.20	38.6%	61.2%
-0.10	38.8%	62.1%
0.00	39.2%	62.7%
+0.10	38.4%	62.3%
+0.20	37.8%	60.9%

636 I Contrastive Steering Directions and the VA Subspace

637 We examine the relationship between contrastive task-specific steering directions and the VA subspace.
 638 The contrastive refusal direction of Arditi et al. [2024] is nearly orthogonal to the VA plane (86.5°).
 639 However, its in-plane component has negative VA coordinates ($V=-0.058$, $A=-0.021$), consistent
 640 with lexical mediation as a contributing mechanism within the contrastive direction.

641 On its target task, the contrastive direction is stronger than VA steering: on HarmBench, refusal
 642 drops to 2.0% at $\alpha = -0.30$ (vs. 59% for VA arousal at $\alpha = +0.30$). However, it causes severe
 643 over-refusal on safe-query benchmarks (OKTest: 19.7% \rightarrow 93.3%, XSTest: 8.4% \rightarrow 94.4% at
 644 $\alpha = +0.45$). When transferred to sycophancy, it produces U-shaped effects—both $+\alpha$ and $-\alpha$
 645 decrease sycophancy (77.6% \rightarrow 57.0% and 77.6% \rightarrow 66.4% respectively)—rather than bidirectional
 646 control.

647 These results illustrate a complementary relationship: task-specific contrastive directions optimize
 648 for single-behavior control, while the VA subspace exposes shared affective structure that enables
 649 monotonic, bidirectional control across behaviors from a single subspace. The near-orthogonality
 650 indicates the model may have multiple possible pathways for modulating one behavior, and we
 651 illuminate how VA subspaces can be one such pathway.

652 J Compute Resources

653 All experiments are conducted on a single node with $8 \times$ NVIDIA H200 GPUs (141 GB HBM3e each).
 654 The compute budget is dominated by the steering sweeps, which evaluate 12 angular directions \times 45
 655 strengths across multiple benchmarks and three models (Llama 3.1-8B, Qwen3-8B, Qwen3-14B).
 656 All individual experiments—including emotion-vector extraction over GoEmotions, per-layer VA
 657 subspace fitting, lexicon projection over 44,728 NRC-VAD words, the full refusal and sycophancy
 658 steering sweeps, logit-lens analyses, neuron ablations, and capability evaluations on MATH-500
 659 and IFEval—finish within a reasonable amount of wall-clock time on this hardware, with no single
 660 experiment exceeding what fits comfortably into the available compute envelope.