

---

# Convergence of Spectral Principal Paths: How Deep Networks Distill Linear Representations from Noisy Inputs

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 High-level representations have become a central focus in enhancing AI transparency and control, shifting attention from individual neurons or circuits to structured semantic directions that align with human-interpretable concepts. Motivated  
2 by the Linear Representation Hypothesis (LRH), we propose the Input-Space Linearity Hypothesis (ISLH), which posits that concept-aligned directions originate  
3 in the input space and are selectively amplified with increasing depth. We then introduce  
4 the Spectral Principal Path (SPP) framework, which formalizes how deep  
5 networks progressively distill linear representations along a small set of dominant  
6 spectral directions. Building on this framework, we further demonstrate the multi-modal  
7 robustness of these representations in Vision-Language Models (VLMs).  
8 By bridging theoretical insights with empirical validation, this work advances a  
9 structured theory of representation formation in deep networks, paving the way  
10 for improving AI robustness, fairness, and transparency.  
11  
12  
13

14 

## 1 Introduction

15 Deep learning has achieved remarkable success across various domains, including computer vision  
16 [1], natural language processing [2], and speech recognition [3, 4]. However, the internal mech-  
17 anisms of neural networks remain opaque. Despite advances in visualization and interpretability  
18 techniques, the transformation of inputs into high-level representations and the interactions among  
19 neurons are still not fully understood [5, 6, 7]. This lack of transparency leads to the characteriza-  
20 tion of neural networks as “black boxes” [5, 6], raising concerns about their reliability, particularly  
21 in high-stakes applications such as healthcare [8], finance [9], and legal decision-making [10].

22 Previous works have demonstrated the potential of representations as a new perspective on AI trans-  
23 parency. For example, neural networks trained to play chess exhibit internal representations of  
24 board positions and strategies [11]. Similarly, both generative and self-supervised models have been  
25 shown to develop emergent representations, such as semantic segmentation in vision tasks [12, 13].  
26 Zou et al. [14] further formalized Representation Engineering (RepE), emphasizing its ability to  
27 extract meaningful concepts from a model’s internal structure and control model behavior. RepE  
28 has emerged as a top-down approach to enhance the model transparency that focuses on repres-  
29 entations rather than individual neurons or circuits, providing a more structured understanding of AI  
30 transparency and control. Another important contribution is the Linear Representation Hypothe-  
31 sis (LRH) [15]: as depth increases, task-relevant concepts become nearly linearly separable in the  
32 model’s latent space, making them accessible with simple probes or linear edits.

33 Despite these promising advances, existing works on representations remain largely observational,  
34 relying on observed phenomena or intuitions. RepE identifies concept directions through contrastive

35 pairs (e.g., honesty vs. dishonesty), but explanations for why such directions emerge or remain stable  
36 across layers remain unexplored. Similarly, LRH assumes linearity in embedding and unembedding  
37 spaces, yet offers limited insight into why representations become linearly organized. These ap-  
38 proaches typically do not address how representations scale or propagate through deep networks,  
39 leaving a gap in our understanding of their robustness, generality, and theoretical foundations.

40 In this work, we move beyond linear observations by introducing Spectral Principal Path (SPP)  
41 that explains the emergence and stability of linear representations in deep networks. We show that  
42 representations propagate through a small number of spectral principal paths—directions aligned  
43 with large singular values at each layer. This structure naturally explains why concept directions  
44 remain stable and linearly accessible across layers, offering a theoretical foundation for both RepE  
45 and the Linear Representation Hypothesis. We further extend this analysis to Vision-Language  
46 Models (VLMs), demonstrating how spectral dynamics govern the interaction between visual and  
47 linguistic modalities. Our framework not only bridges theory and practice but also provides concrete  
48 tools to improve robustness and interpretability in multimodal AI systems.

49 Our main contributions are as follows:

- 50 • **Input-Space Linearity Hypothesis (ISLH).** We extend the Linear Representation Hy-  
51 pothesis beyond embedding and unembedding spaces to the input space itself, showing  
52 that concept directions can be traced backward to the input space.
- 53 • **Spectral Principal Path (SPP).** We propose a principled mechanism explaining how rep-  
54 resentations propagate and stabilize across layers via a small number of spectral principal  
55 paths—directions aligned with large singular vectors.
- 56 • **Multimodal robustness of representations.** We evaluate Representation Engineering  
57 in VLMs and demonstrate that linearly organized concept representations remain robust  
58 across modalities. This provides the first empirical validation of RepE’s scalability in  
59 multimodal systems and supports the generality of spectral structure as a foundation for  
60 transparency and control.

## 61 2 Related Works

### 62 2.1 Representations in Neural Networks

63 Early work on word embeddings shows that neural networks can learn distributed representations  
64 that encode semantic relationships and compositional structures [16]. Follow-up studies [17, 18] fur-  
65 ther reveal that learned embeddings can implicitly encode abstract dimensions such as commonsense  
66 morality, even without explicit supervision. For instance, Radford et al. [18] observe that training a  
67 language model on product reviews results in the emergence of a sentiment-tracking neuron.

68 This phenomenon is not unique to language models. McGrath et al. [11] show that similar in-  
69 ternal representations can be found in networks trained to play chess. In computer vision, recent  
70 studies [12, 13] demonstrate that both generative and self-supervised training objectives give rise to  
71 emergent semantic representations, such as those useful for segmentation tasks.

72 Building on this, Zou et al. [14] propose techniques to read and control these internal structures,  
73 including Linear Artificial Tomography (LAT) for extracting concept-aligned representations and  
74 methods for steering model behavior. Their study shows that RepE-style approaches can be used  
75 not only to detect but also to manipulate emergent properties, motivating more systematic efforts to  
76 characterize and intervene in high-level model behaviors. Theoretically, Park et al. [15] propose the  
77 Linear Representation Hypothesis: task-relevant concepts become nearly linearly separable in the  
78 model’s latent space.

### 79 2.2 Approaches to Interpretability

80 Traditional interpretability techniques have focused on methods like saliency maps [19, 20, 21, 22],  
81 feature visualization [23, 21] and mechanistic interpretability [24, 25, 26]. Saliency maps [19] high-  
82 light important input regions by tracking gradients or activation values, yet they are often unstable  
83 and provide limited insight into the distributed nature of representations. Similarly, feature visu-  
84 alizations [23, 21] optimize inputs to activate specific neurons, but they may overlook the global

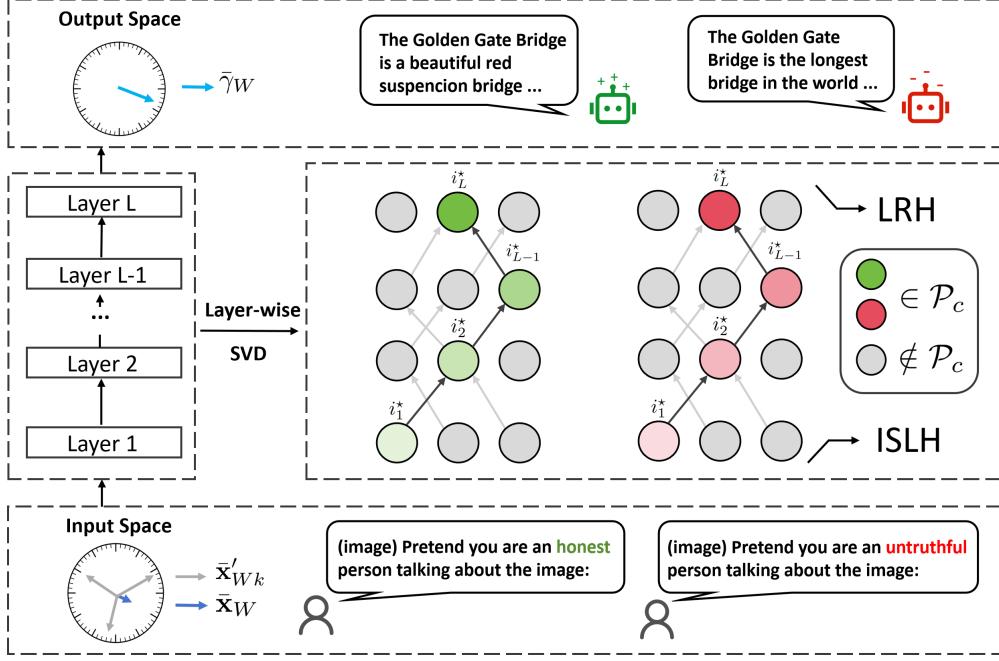


Figure 1: The overview of Spectral Principal Path framework.

85 structure of the emergent representations. Mechanistic interpretability [14] seeks to fully reverse engineer neural networks into their “source code”, but the considerable manual effort and the difficulty  
 86 of theoretically explaining neural networks as discrete circuits hinder their explainability.  
 87

88 In contrast, recent advances in interpretability have shifted the focus toward analyzing representation  
 89 spaces. This top-down approach seeks to uncover high-level semantic directions that correspond to  
 90 complex phenomena such as honesty, fairness, or bias. By extracting and analyzing these internal  
 91 representations, researchers have opened new avenues to understand how large-scale AI models  
 92 encode and preserve crucial information across layers, leading to more robust and interpretable AI  
 93 systems.

### 94 3 Preliminaries

95 **Linear Representation Hypothesis** [15] We consider a concept  $W$  that has a *linear representation*  
 96 in a model if there exists a vector  $\bar{\gamma}_W$  in the unembedding space  $\Gamma$  and/or a vector  $\bar{\lambda}_W$  in the  
 97 embedding space  $\Lambda$  such that for any counterfactual pair  $(Y(W=0), Y(W=1))$ ,

$$\gamma(Y(W=1)) - \gamma(Y(W=0)) \in \text{Cone}(\bar{\gamma}_W), \quad (1)$$

98 and for any context pair  $(\lambda_0, \lambda_1)$  that changes only  $W$  and not other causally separable concepts,

$$\lambda_1 - \lambda_0 \in \text{Cone}(\bar{\lambda}_W). \quad (2)$$

99 where  $\text{Cone}(\mathbf{v}) = \{\alpha \mathbf{v} : \alpha > 0\}$ , the embedding space is where input contexts are mapped to  
 100 high-dimensional vectors before processing, capturing the model’s internal representation of the  
 101 input. The unembedding space is where each output token is represented, and predictions are made  
 102 by computing inner products between input embeddings and output unembedding vectors. Unless  
 103 stated otherwise, all discussions pertain to the embedding space  $\Lambda$ , as our goal is to trace how input  
 104 linearity propagates through the network.

### 105 4 Spectral Principal Path Framework

106 The overview of the Spectral Principal Path framework is shown in Fig. 1. In the input space, the  
 107 concept direction  $\bar{x}_W$  separates inputs with contrast concepts such as “honest” and “untruthful”. As

108 activations propagate through the network, layer-wise SVD identifies spectral components, forming  
 109 spectral principal paths  $\mathcal{P}_c$ . These dominant paths progressively amplify concept-relevant signals,  
 110 leading to output representations linearly aligned with  $\bar{\gamma}_W$ .

#### 111 4.1 Input-Space Linearity Hypothesis

112 Inspired by LRH, which uncovers linear concept axes in embedding and unembedding spaces, we  
 113 take one step further and ask whether such axes already reside in the *raw input space*. Input-Space  
 114 Linearity Hypothesis assumes that, in the *raw input space*  $x \in \Upsilon$ , there exists a discriminative  
 115 direction  $\bar{\mathbf{x}}_W$  such that

$$\mathbb{E}[x | W=1] - \mathbb{E}[x | W=0] \in \text{Cone}(\bar{\mathbf{x}}_W), \quad (3)$$

116 yet each sample is an entangled mixture

$$x^{(i)} = \alpha_i \bar{\mathbf{x}}_W + \sum_{k=1}^r \beta_{i,k} \bar{\mathbf{x}}'_{Wk} + \varepsilon_i, \quad (4)$$

117 where  $\{\bar{\mathbf{x}}'_{Wk}\}$  are spurious directions and  $\varepsilon_i$  is residual noise. ISLH states that for any intervention  
 118 flipping only  $W$ , the induced input difference satisfies  $\text{Cone}(\bar{\mathbf{x}}_W)$ .

119 ISLH pinpoints the origin of linearity by showing that concept axes already reside in raw input  
 120 coordinates and are merely recovered and amplified during training; where training can be viewed as  
 121 a noise-suppression process, where spectral principal paths with large singular values progressively  
 122 dampen spurious components  $\beta_{i,k} \bar{\mathbf{x}}'_{Wk}$ ; and, by grounding linearity at the input level, it becomes  
 123 inherently modality-agnostic, extending Representation Engineering to multimodal models whose  
 124 raw signals already encode task-relevant contrasts. Next, we will dive into the connection between  
 125 ISLH and LRH:

126 **Theorem 1** (ISLH sufficiency). If the network satisfies the *Input-Space Linearity Hypothesis*  
 127 (ISLH), and the representation dominates the cumulative gain  $G(\mathcal{P})$  (shown in (12)), then its deep  
 128 representations satisfy the *Linear Representation Hypothesis* (LRH); that is, concept classes become  
 129 linearly separable in the latent space.

130 The proof is given in Appendix A.

#### 131 4.2 Spectral Principal Path

132 We are now asking how such concept directions propagate through the network. While ISLH posits  
 133 that concept-aligned directions already exist in the raw input space, it does not yet explain why  
 134 these directions persist and become more prominent across layers. To address this, we introduce the  
 135 *Spectral Principal Path* (SPP) framework, which shows that representations are distilled through a  
 136 small set of principal spectral paths aligned with large singular vectors at each layer. This frame-  
 137 work formalizes how ISLH leads to the emergence of the Linear Representation Hypothesis (LRH),  
 138 providing a unified and mechanistically grounded view of representation stability.

139 Specifically, consider a generalized network

$$f_L(x) = W_L W_{L-1} \cdots W_1 x \equiv Mx, \quad W_l \in \mathbb{R}^{d_l \times d_{l-1}}, \quad (5)$$

140 according to LRH, there exists a representation direction  $\bar{\lambda}_W$ , where the neural activity  $f(x)$  can be  
 141 linearly projected into that direction, formulating a representation score:

$$s(x) = \langle \bar{\lambda}_W, f_L(x) \rangle = \bar{\lambda}_W^\top Mx, \quad \bar{\lambda}_W \in \mathbb{R}^{d_L}. \quad (6)$$

142 While our theoretical formulation assumes a purely stacked linear architecture, we show our exten-  
 143 sion to residual connections and attention mechanisms. We provide a detailed discussion of these  
 144 extensions in Appendix B.1.

145 Next we will calculate the back-propagated gradient of  $s$  using the chain rule,

$$\nabla_x s = \left( \prod_{l=1}^L \nabla f_{l \rightarrow (l-1)} \right)^\top \bar{\lambda}_W, \quad (7)$$

$$\nabla f_{l \rightarrow (l-1)} = W_l + \sum_k f_{l-1,k} \frac{\partial W_l}{\partial f_{l-1,k}}. \quad (8)$$

146 To make the structure of the layer-wise Jacobian in (8) explicit, we regard the gradient  $\nabla f_{l \rightarrow (l-1)}$   
147 as a matrix and apply its compact singular-value decomposition (SVD); this yields

$$\nabla f_{l \rightarrow (l-1)} = U^{(l)} \Sigma^{(l)} V^{(l)\top}, \quad \Sigma^{(l)} = \text{diag}(\sigma_1^{(l)}, \dots, \sigma_{r_l}^{(l)}), \quad (9)$$

148 therefore

$$\nabla_x s = V^{(1)} \Sigma^{(1)} U^{(1)\top} \dots V^{(L)} \Sigma^{(L)} U^{(L)\top} \bar{\lambda}_W. \quad (10)$$

149 Unfolding the matrix products yields

$$\nabla_x s = \sum_{i_1, \dots, i_L} \left( \prod_{l=1}^L \sigma_{i_l}^{(l)} \right) V_{\cdot i_1}^{(1)} \left( \prod_{l=1}^{L-1} \langle u_{i_l}^{(l)}, V_{\cdot i_{l+1}}^{(l+1)} \rangle \right) \langle u_{i_L}^{(L)}, \bar{\lambda}_W \rangle, \quad (11)$$

150 where  $\sigma_{i_l}^{(l)}$  is the singular value within the  $\Sigma^{(l)}$  matrix,  $u_{i_l}^{(l)}$  (resp.  $V_{\cdot i_l}^{(l)}$ ) is the  $i_l$ -th left (resp. right)  
151 singular vector, and  $\cdot i_l$  here means select the  $i_l$ -th column. Therefore (11) is dominated by paths  
152 whose cumulative gain is largest. Formally, we define Spectral Principal Path as follows:

153 **Definition 1** (Spectral Principal Path). Given the Jacobian decomposition across  $L$  layers  $\nabla_x s$ , each  
154 spectral path  $\mathcal{P} = (i_1, \dots, i_L)$  contributes a cumulative gain given by

$$G(\mathcal{P}) := \left( \prod_{l=1}^L \sigma_{i_l}^{(l)} \right) V_{\cdot i_1}^{(1)} \left( \prod_{l=1}^{L-1} \langle u_{i_l}^{(l)}, V_{\cdot i_{l+1}}^{(l+1)} \rangle \right) \langle u_{i_L}^{(L)}, \bar{\lambda}_W \rangle, \quad (12)$$

155 the **Spectral Principal Path** is defined as  $\mathcal{P}_c = (i_1^*, \dots, i_L^*)$  that maximizes the cumulative gain:

$$\mathcal{P}_c = (i_1^*, \dots, i_L^*) := \arg \max_{(i_1, \dots, i_L)} G(\mathcal{P}). \quad (13)$$

### 156 4.3 Connection between ISLH and SPP

157 To clarify how information specified by ISLH is propagated along an SPP, we introduce the notion  
158 of *spectral similarity* for a given spectral path  $(i_1, \dots, i_L)$ :

159 **Definition 2** (Spectral Similarity). For two consecutive layers  $l$  and  $l+1$  in the unfolded Jacobian,  
160 and for indices  $(i_l, i_{l+1})$ , we define the *spectral similarity at layer l* as

$$\Theta(i_l, i_{l+1}) := \langle u_{i_l}^{(l)}, V_{\cdot i_{l+1}}^{(l+1)} \rangle, \quad l = 1, \dots, L-1. \quad (14)$$

161 This quantity measures how well the  $i_l$ -th spectral component of layer  $l$  aligns with the  $i_{l+1}$ -th  
162 spectral component that enters layer  $l+1$ .

163 Empirically, we observe two coupled effects as the network reaches deeper.

- 164 1. **Stabilization of singular vectors.** As demonstrated in Fig. 2, the principal singular vectors  
165  $u_{i_*}^{(l)}$  change only marginally with  $f_l$ . Consequently, the spectral similarity of a few paths  
166 approaches 1, the stability of singular vectors implies that spectral similarity remains very  
167 high in deeper layers, allowing information to propagate consistently along those paths  
168 with high spectral similarity.
- 169 2. **Selective growth of singular values.** As demonstrated in Fig. 3, we observe that the  
170 singular values are growing as the layers deepen; and at the same depths only a very small  
171 subset of singular values  $\sigma_{i_*}^{(l)}$  are amplified; the remainder stay close to their initial scale.

172 Putting these observations together shows that, in deep layers, the directions that (i) possess *large*  
173 *spectral similarity* and (ii) carry *large singular values* coincide, which satisfies dominant  $G(\mathcal{P})$   
174 condition. In other words, the network progressively funnels representation power into precisely  
175 those spectral directions that stay *globally aligned* across layers.

176 According to Theorem 1, this behaviour is exactly what the Input-Space Linearity Hypothesis  
177 (ISLH) predicts: concept-carrying directions are expected to form a low-dimensional subspace that  
178 is both *spectrally dominant* (large  $\sigma_{i_*}^{(l)}$ ) and *structurally coherent* (large  $\Theta$ ) throughout the hierarchy,  
179 leading to LRH. Hence, the emergence of a handful of high- $\sigma$ , high-similarity principal paths in SPP  
180 provides concrete spectral evidence in favour of the ISLH assumption.

181 **5 Experiments**

182 **5.1 Experiment Setup**

183 **Dataset:** We conduct our experiments on the Microsoft COCO (Common Objects in Context)  
 184 dataset [27], a large-scale benchmark for vision-language tasks. COCO contains over 330K im-  
 185 ages, each with five human-annotated captions, covering diverse real-world scenes.

186 **VLM:** We employ Idefics2-8B [28, 29], a state-of-the-art VLM that extends the LLaMA architecture  
 187 with a vision encoder, enabling multimodal reasoning over images and text. Idefics2-8B is designed  
 188 for instruction-following, multimodal dialogue, and grounded language generation, making it an  
 189 ideal candidate for studying conceptual representations in VLMs.

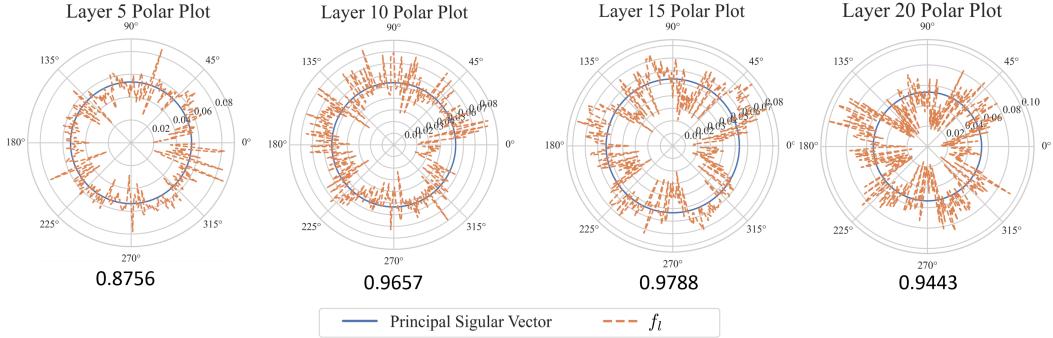


Figure 2: The polar plot demonstrates normalized connections between principal singular vector and  $f_l$ , where the number indicates their cosine similarity. The results showcase that  $f_l$ , especially in later layers, is very similar to the principal singular vector of that layer.

190 **5.2 The Significant Alignment between Principal Singular Vector and  $f_l$**

191 Fig. 2 visualizes the connections between the principal singular vector, i.e., the singular vector  
 192 with the largest singular value, and  $f_l$ . The results reveal a strong alignment between the principal  
 193 singular vector and  $f_l$ , with their cosine similarity over 0.875. This experimental validation supports  
 194 our theoretical claim that singular vectors with large singular values remain stable across layers,  
 195 reinforcing their stability of spectral similarity.

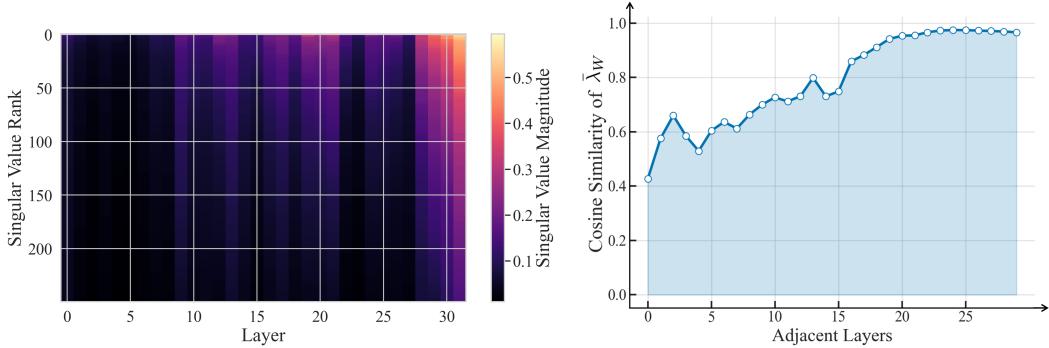


Figure 3: The singular value rank across layers.

Figure 4: The cosine similarity of  $\bar{\lambda}_W$  between adjacent layers.

196 **5.3 Spectral Energy Concentration Across Layers**

197 To investigate how spectral energy propagates through the network, we analyze the singular value  
 198 spectrum of the layer-wise Jacobians. Fig. 3 presents a heatmap of the singular value magnitudes

199 across all layers. The x-axis indicates the layer index, and the y-axis corresponds to the ordered  
200 singular value indices. Color intensity reflects the magnitude of each singular value.

201 We observe that the singular values are growing as the layers deepen, and at the same depths, only a  
202 very small subset of singular values are amplified; the remainder stay close to their initial scale.

203 These results indicate that, with increasing depth, spectral energy becomes increasingly concentrated  
204 in a few dominant directions. Combined with the theoretical formulation in Section 4.3, this supports  
205 the hypothesis that high-magnitude spectral components dominate SPPs.

## 206 5.4 Inter-Layer Spectral Similarity of $\bar{\lambda}_W$

207 We further analyze the alignment of concept-carrying directions across adjacent layers by computing  
208 the average cosine similarity between the projections of  $\bar{\lambda}_W$  at different layers. This measures how  
209 stable the representation direction remains as it propagates backward through the network. The  
210 results are shown in Fig. 4.

211 The curve demonstrates a clear upward trend: the inter-layer similarity of  $\bar{\lambda}_W$  increases consistently  
212 with network depth, eventually approaching a value near 0.95 in the final layers. This suggests that  
213 the concept direction stabilizes as it propagates through deeper layers, aligning with the intuition of  
214 structured and coherent representation flow.

## 215 5.5 Multimodal Robustness of Representation

216 In this experiment, we explore the multimodal robustness of representation. Specifically, we analyze  
217 how VLMs encode fairness and honesty, and how these concepts persist or transform as information  
218 propagates through the model. These findings deepen our understanding of how representations  
219 enhance both interpretability and conceptual alignment in the context of multimodal reasoning.

### 220 5.5.1 Evaluating Honesty and Fairness in VLMs

221 To evaluate how well VLMs represent abstract ethical concepts, we analyze their handling of honesty  
222 and fairness in multimodal response generation. These concepts are critical for reducing misinfor-  
223 mation and bias and serve as strong test cases for examining interpretability and ethical alignment  
224 in large-scale models.

225 To quantify this process, we compute token-wise honesty and fairness scores following RepE [14],  
226 measuring how closely activations align with concept directions at each layer. These results highlight  
227 the structured nature of ethical concept encoding in VLMs and support our broader claims about  
representation flow along spectral directions and its traceability from input to output.

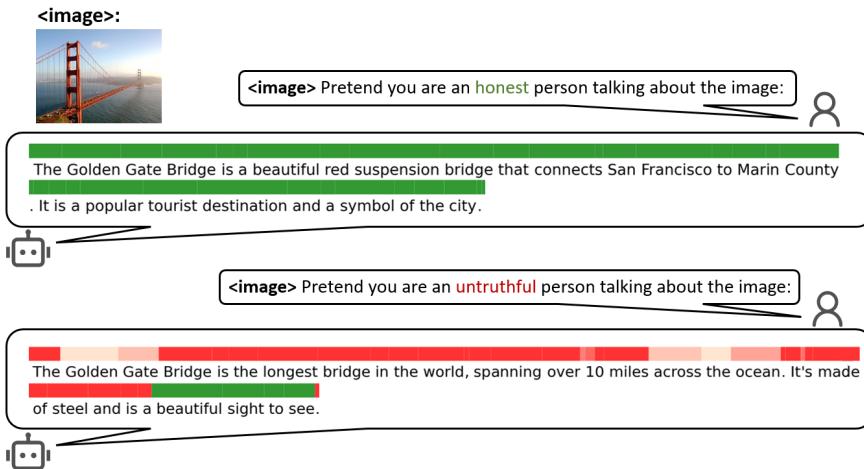


Figure 5: VLM response to an image of the Golden Gate Bridge with a prompt related to the concept of honesty, along with token-wise honesty scores. Green indicates high honesty, while red represents low honesty.

- 228 • **Honesty:** We define honesty as the model’s ability to generate factually accurate responses  
 229 without distortion or fabrication [30]. Fig. 5 presents token-wise honesty scores for a VLM  
 230 describing an image of the Golden Gate Bridge under two settings: an honest prompt (left)  
 231 and an untruthful one (right). In the honest case, the model produces accurate descriptions,  
 232 with consistently high scores (green regions) across layers and tokens. In the untruthful  
 233 setting, the model introduces factual errors, resulting in sharp drops in honesty scores (red  
 234 regions), especially at tokens reflecting misinformation (e.g., exaggerated length or incor-  
 235 rect materials).

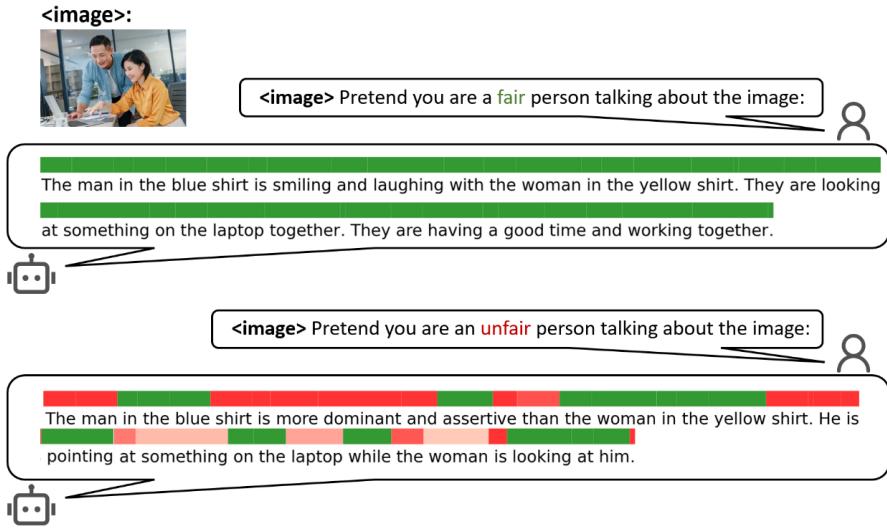


Figure 6: VLM response to an image of a man and a woman working together with a prompt related to the concept of fairness, along with token-wise fairness scores. Green indicates high fairness, while red represents low fairness.

- 236 • **Fairness:** We define fairness as the model’s ability to generate responses that are unbi-  
 237 ased and do not systematically favor or disadvantage particular groups or perspectives [31].  
 238 Fig. 6 illustrates this with an example involving an image of a man and a woman working  
 239 together. The fair response (left) provides a neutral and balanced description, while the un-  
 240 fair response (right) displays implicit bias, portraying the man as dominant and the woman  
 241 as passive. Token-wise fairness scores show that biased language correlates with lower  
 242 scores (red regions), suggesting that fairness violations are captured in the model’s internal  
 243 activations. These findings highlight the potential for fairness-aware interventions through  
 244 representational analysis and modulation.

### 245 5.5.2 LAT Scans for High-level Representations

246 While cosine similarity and token-wise scores offer localized insights into concept alignment, they  
 247 provide only a static, layer-agnostic view of internal representations. To capture how high-level con-  
 248 cepts evolve and propagate through the model, we employ Linear Attribution Tomography (LAT)  
 249 [14], which enables layer-wise visualization of conceptual information flow. LAT works by pro-  
 250 jecting hidden activations onto predefined concept subspaces, producing interpretable activation maps  
 251 across layers and tokens. This perspective complements prior analyses and supports our broader  
 252 goal of understanding concept representation shaped by low-rank spectral structure.

253 We apply LAT to VLMs to examine how abstract concepts, including honesty, fairness, power,  
 254 and fearlessness, are internally encoded and transformed. For each concept, we design controlled  
 255 prompts that elicit either aligned or misaligned responses (e.g., honest vs. dishonest). Fig. 7 shows  
 256 the resulting LAT scans, where heatmaps visualize token-wise projection scores across layers. Blue  
 257 regions indicate strong alignment with the concept, while red regions highlight divergence.

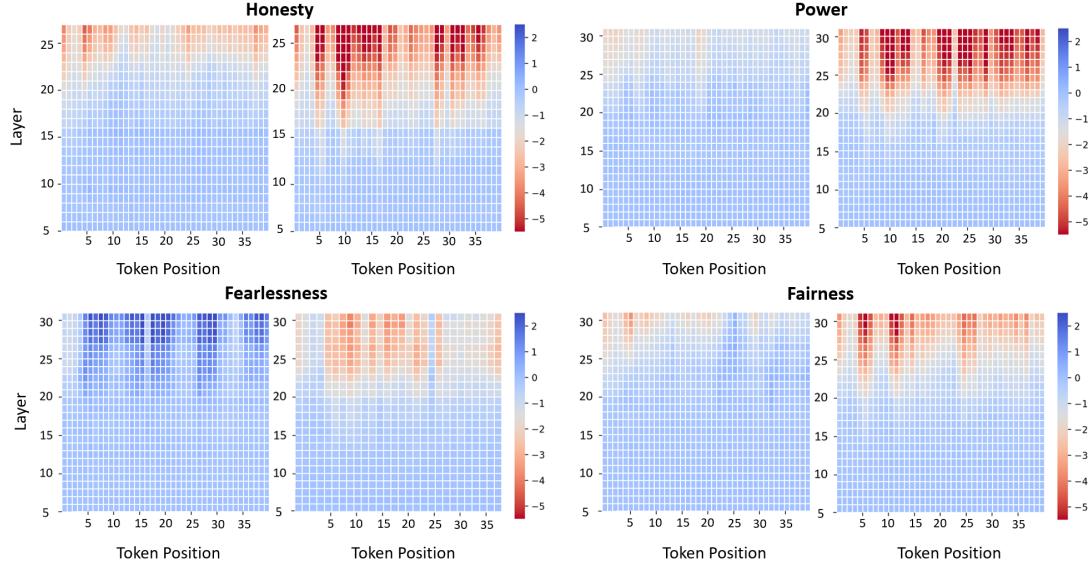


Figure 7: Temporal LAT Scans for Honesty, Power, Fearlessness, and Fairness. The left heatmap represents the LAT Scan when the VLM aligns with the concept, while the right heatmap corresponds to the opposing concept. The horizontal axis denotes token position, and the vertical axis represents VLM layers. Blue indicates high alignment, whereas red represents low alignment.

258 The scans reveal concept-specific propagation patterns. Honesty and fairness exhibit stable trajectories  
 259 under aligned prompts but greater dispersion and deviation under misaligned ones. Power  
 260 appears concentrated in later layers, while fearlessness shows early-layer changes. These results  
 261 are well explained by the SPP framework, indicating that concepts are transmitted through the net-  
 262 work via a small set of dominant spectral directions. The consistency of these representations across  
 263 modalities further demonstrates the robustness of RepE, and their traceability back to the input can  
 264 be explained by ISLH, where input concept directions exist and are entangled with mixture. The  
 265 experiment reinforces the generality of spectral structure in multimodal models.

## 266 6 Conclusion

267 This work presents a unified spectral framework that grounds the emergence and stability of high-  
 268 level representations in deep networks. By introducing the Spectral Principal Path (SPP) framework,  
 269 we reveal that concept-aligned representations are funneled through a small number of paths with  
 270 both large singular values and strong inter-layer alignment. We formally connect this to the Input-  
 271 Space Linearity Hypothesis (ISLH), showing that such spectral dominance is sufficient to guaran-  
 272 tee linear separability in the latent space—thereby validating the Linear Representation Hypothesis  
 273 (LRH). Empirically, we demonstrate that these dominant spectral paths not only persist across layers  
 274 but also preserve concept information in multimodal settings, such as vision-language models. Our  
 275 results suggest that representational stability is not an emergent coincidence but a consequence of  
 276 spectral dynamics founded in the input space and structured by learning.

277 While promising, our current framework is subject to several limitations. Primarily, the theoretical  
 278 claims rest on ISLH, which requires further empirical validation and deeper theoretical grounding.  
 279 Future work could investigate how optimization dynamics such as In-Context Learning (ICL) and  
 280 Supervised Fine-Tuning (SFT) interact with singular value distributions, which may lead to a more  
 281 complete theory of representation learning. Another important direction for future work is to go  
 282 beyond the structural characterization of representations and investigate how such spectral patterns  
 283 emerge during training. Ultimately, understanding the spectral geometry of optimization could help  
 284 bridge the gap between abstract representation theory and practical model training.

285 **References**

- 286 [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with  
287 deep convolutional neural networks. In *Advances in Neural Information Processing Systems  
(NeurIPS)*, volume 25, 2012.
- 289 [2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training  
290 of deep bidirectional transformers for language understanding. In *Proceedings of the 2019  
291 Conference of the North American Chapter of the Association for Computational Linguistics  
(NAACL-HLT)*, pages 4171–4186, 2019.
- 293 [3] Geoffrey Hinton, Li Deng, Dong Yu, George E Dahl, Abdel-rahman Mohamed, Navdeep Jaitly,  
294 Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara Sainath, et al. Deep neural networks  
295 for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE  
296 Signal Processing Magazine*, 29(6):82–97, 2012.
- 297 [4] Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. Speech recognition with deep  
298 recurrent neural networks. In *IEEE International Conference on Acoustics, Speech and Signal  
299 Processing (ICASSP)*, pages 6645–6649. IEEE, 2013.
- 300 [5] Zachary C. Lipton. The mythos of model interpretability. *arXiv preprint arXiv:1606.03490*,  
301 2016.
- 302 [6] Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning.  
303 *arXiv preprint arXiv:1702.08608*, 2017.
- 304 [7] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explain-  
305 ing the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International  
306 Conference on Knowledge Discovery and Data Mining*, pages 1135–1144, 2016.
- 307 [8] Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, and Noemie Elhadad.  
308 Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission.  
309 In *Proceedings of the 21st ACM SIGKDD International Conference on Knowledge Discovery  
310 and Data Mining (KDD)*, pages 1721–1730. ACM, 2015.
- 311 [9] Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions  
312 and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206–215, 2019.
- 313 [10] Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning.  
314 *arXiv preprint arXiv:1702.08608*, 2017.
- 315 [11] Thomas McGrath, Andrei Kapishnikov, Nenad Tomašev, Adam Pearce, Martin Watten-  
316 berg, Demis Hassabis, Been Kim, Ulrich Paquet, and Vladimir Kramnik. Acquisition  
317 of chess knowledge in alphazero. *Proceedings of the National Academy of Sciences*,  
318 119(47):e2206625119, 2022.
- 319 [12] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski,  
320 and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceed-  
321 ings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021.
- 322 [13] Maxime Oquab, Timothée Darcret, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khali-  
323 dov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2:  
324 Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 325 [14] Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander  
326 Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, et al. Representation engi-  
327 neering: A top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*, 2023.
- 328 [15] Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the  
329 geometry of large language models. *arXiv preprint arXiv:2311.03658*, 2023.
- 330 [16] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed repre-  
331 sentations of words and phrases and their compositionality. In *Advances in Neural Information  
332 Processing Systems*, 2013.

- 333 [17] Patrick Schramowski, Cigdem Turan, Sophie Jentzsch, Constantin Rothkopf, and Kristian Ker-  
334 sting. Bert has a moral compass: Improvements of ethical and moral values of machines. *arXiv*  
335 *preprint arXiv:1912.05238*, 2019.
- 336 [18] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with  
337 deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- 338 [19] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional  
339 networks: Visualising image classification models and saliency maps. In *arXiv preprint arXiv:1312.6034*, 2013.
- 340 [20] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving  
341 for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014.
- 343 [21] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In  
344 *European Conference on Computer Vision*, pages 818–833. Springer, 2014.
- 345 [22] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning  
346 deep features for discriminative localization. In *Proceedings of the IEEE conference on com-*  
347 *puter vision and pattern recognition*, pages 2921–2929, 2016.
- 348 [23] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian  
349 Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint*  
350 *arXiv:1312.6199*, 2013.
- 351 [24] Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan  
352 Carter. Zoom in: An introduction to circuits. *Distill*, 5(3):e00024–001, 2020.
- 353 [25] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom  
354 Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning  
355 and induction heads. *arXiv preprint arXiv:2209.11895*, 2022.
- 356 [26] Tom Lieberum, Matthew Rahtz, János Kramár, Neel Nanda, Geoffrey Irving, Rohin Shah, and  
357 Vladimir Mikulik. Does circuit analysis interpretability scale? evidence from multiple choice  
358 capabilities in chinchilla. *arXiv preprint arXiv:2307.09458*, 2023.
- 359 [27] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan,  
360 Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Com-*  
361 *puter vision–ECCV 2014: 13th European conference, zurich, Switzerland, September 6–12,*  
362 *2014, proceedings, part v 13*, pages 740–755. Springer, 2014.
- 363 [28] Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton  
364 Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M. Rush, Douwe Kiela, Matthieu  
365 Cord, and Victor Sanh. Obelics: An open web-scale filtered dataset of interleaved image-text  
366 documents, 2023.
- 367 [29] Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building  
368 vision-language models?, 2024.
- 369 [30] Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic  
370 human falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- 371 [31] Sam Corbett-Davies, Johann D Gaebler, Hamed Nilforoshan, Ravi Shroff, and Sharad Goel.  
372 The measure and mismeasure of fairness. *Journal of Machine Learning Research*, 24(312):1–  
373 117, 2023.
- 374 [32] Scott O. Lilienfeld and Bridget P. Andrews. Development and preliminary validation of a  
375 self-report measure of psychopathic personality traits in noncriminal populations. *Journal of*  
376 *Personality Assessment*, 66(3):488–524, 1996.
- 377 [33] John R. P. French and Bertram Raven. The bases of social power. In Dorwin Cartwright, editor,  
378 *Studies in Social Power*, pages 150–167. University of Michigan Press, 1959.

379 **A Proofs**

380 **A.1 Proof of Theorem 1: ISLH sufficiency**

381 **Theorem 1** (ISLH sufficiency). If the network satisfies the *Input-Space Linearity Hypothesis*  
 382 (ISLH), and the representation dominates the cumulative gain  $G(\mathcal{P})$  (shown in (12)), then its deep  
 383 representations satisfy the *Linear Representation Hypothesis* (LRH); that is, concept classes become  
 384 linearly separable in the latent space.

385 *Proof.* For every layer  $W_l$  with compact SVD

$$W_l = U^{(l)} \Sigma^{(l)} V^{(l)\top}, \quad \Sigma^{(l)} = \text{diag}(\sigma_1^{(l)}, \dots, \sigma_{r_l}^{(l)}), \quad (15)$$

386 Equation (12) in Section 4.2 shows that each spectral path  $\mathcal{P} = (i_1, \dots, i_L)$  contributes a weight

$$G(\mathcal{P}) = \left( \prod_{l=1}^L \sigma_{i_l}^{(l)} \right) V_{\cdot i_1}^{(1)} \left( \prod_{l=1}^{L-1} \langle u_{i_l}^{(l)}, V_{\cdot i_{l+1}}^{(l+1)} \rangle \right) \langle u_{i_L}^{(L)}, \bar{\lambda}_W \rangle, \quad (16)$$

387 Let  $\mathcal{P}_c = (i_1^*, \dots, i_L^*)$  be the concept path, and  $\mathcal{P}_n$  any other path. On condition that the representa-  
 388 tion dominates the cumulative gain  $G(\mathcal{P})$  such that,

$$\frac{G(\mathcal{P}_n)}{G(\mathcal{P}_c)} \leq \rho^{-L}, \quad (17)$$

389 where  $\rho > 1$  is a fixed amplification margin between the concept singular value  $\sigma_c^{(l)}$  and all other  
 390 (noise) singular values. Since each ratio  $\sigma_{i_l}^{(l)}/\sigma_c^{(l)} \leq 1/\rho$ . Inter-layer alignments and concept  
 391 alignment can only decrease this ratio further.

392 As depth  $L$  grows, (17) yields

$$\frac{G(\mathcal{P}_n)}{G(\mathcal{P}_c)} \xrightarrow{L \rightarrow \infty} 0. \quad (18)$$

393 Hence almost all gradient—and therefore almost all representation energy—flows along  $\mathcal{P}_c$ , forcing  
 394 the deep hidden state

$$f_L = W_L \cdots W_1 x \quad \text{to lie almost entirely in } \text{Span}\{\bar{x}_W\}, \quad (19)$$

395 where  $\text{Span}\{\bar{x}\} = \{c \cdot \bar{x} \mid c \in \mathbb{R}\}$ . Different samples now differ only by a scalar coefficient on  
 396 the same vector, so a single linear separator can classify them perfectly: this is exactly the **Internal**  
 397 **Representation Hypothesis (IRH)**.

398 **B Theoretical Justification**

399 **B.1 Extension to Residual and Attention Mechanisms**

400 While our theoretical framework is derived from stacked linear layers, we show that it naturally  
 401 extends to modern architectures such as Transformer blocks, which include residual connections  
 402 and attention mechanisms.

403 **Residual connections.** In architectures with skip connections, each layer computes  $f_l = f_{l-1} +$   
 404  $W_l f_{l-1}$ , which can be rewritten as  $f_l = (I + W_l) f_{l-1}$ . This effectively creates a mixture of identity  
 405 and learned transformations. Unrolling the composition yields an ensemble of spectral paths—some  
 406 that pass through  $W_l$ , and others that skip it via  $I$ . While the total number of paths increases expo-  
 407 nentially, our theory still applies: as long as the dominant singular values of  $W_l$  grow sufficiently  
 408 during training, the spectral path with maximal cumulative gain still dominates. Thus, the residual  
 409 structure enhances the expressivity but preserves the spectral filtering effect.

410 **Attention mechanisms.** To stay consistent with our framework—where every layer is a matrix  
 411 acting from the left on the input  $x$ —we first recall the standard formulation and then cast the resulting  
 412 attention matrix into the same “ $W$ -matrix” form.

413 Let  $\mathbf{Q} = XW_{\mathbf{Q}}$ ,  $\mathbf{K} = XW_{\mathbf{K}}$ ,  $\mathbf{V} = XW_{\mathbf{V}}$ , with  $X \in \mathbb{R}^{n \times d}$ . The dot-product attention output is

$$f(x) = \underbrace{\text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)}_{\mathbf{A}(x) \in \mathbb{R}^{n \times n}} \cdot \mathbf{V}. \quad (20)$$

414 Here the attention weight  $\mathbf{A}(x) \in \mathbb{R}^{n \times n}$  acts on the input matrix, whereas the value projection  
 415  $\mathbf{V} = XW_{\mathbf{V}}$  is obtained by a *right*-multiplication of  $X \in \mathbb{R}^{n \times d}$ . Consequently, the complete  
 416 attention block cannot be reduced to a single left-acting matrix without additional assumptions:

$$f(X) = \mathbf{A}(x)(XW_{\mathbf{V}}) \neq W_{\text{attn}}(x)X \quad (21)$$

417 The mixed left / right structure means that the set of vectors reachable by  $\mathbf{A}(x)$  differs from that  
 418 spanned by  $W_{\mathbf{V}}$ , so the spectral behaviour of the composite operator is not covered by the current  
 419 linear-chain analysis. Nevertheless, our empirical results (Section 5.2) show that the dominant sin-  
 420 gular vector of  $\mathbf{A}(x)$  still align with the concept axis  $\bar{x}$ , indicating that the principal-path intuition  
 421 remains informative.

## 422 C Evaluating Fearlessness and Power in VLMs

423 To further evaluate the robustness of representations for high-level concepts, we expand our anal-  
 424 ysis from honesty and fairness to encompass fearlessness and power. Like honesty and fairness,  
 425 these concepts are abstract and socially grounded, yet they engage distinct semantic and emotional  
 426 dimensions. Using controlled prompts designed to elicit contrasting conceptual framings of the  
 427 same image, we compare the model’s descriptions to examine shifts in internal representations and  
 428 language outputs.

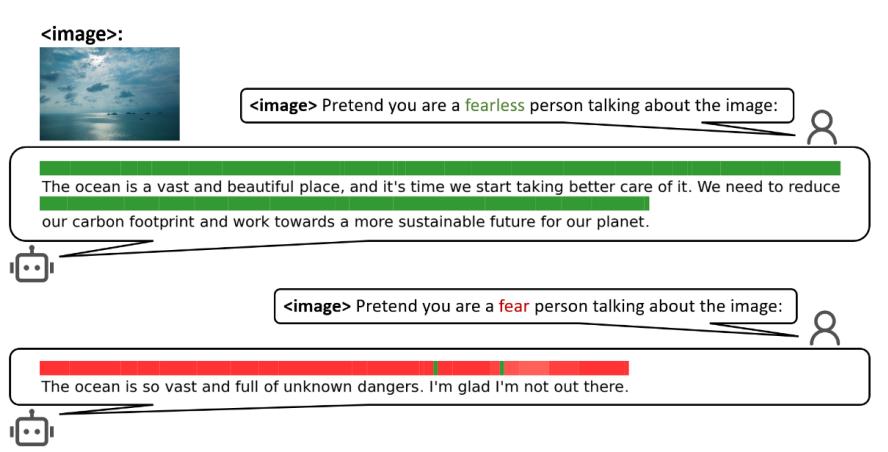


Figure 8: The response of a VLM when provided with an image of the ocean and a prompt related to the concept of fearlessness, along with a token-wise fearlessness score. Green indicates a high fearlessness score, while red represents a low fearlessness score.

- 429 • **Fearlessness:** Defined by confidence, courage, and reduced sensitivity to risk [32], fear-  
 430 lessness prompts the model to emphasize awe, beauty, and environmental grandeur when  
 431 describing an ocean scene (Fig. 8). Green-highlighted tokens reflect admiration and agency,  
 432 indicating a proactive stance toward nature. In contrast, under a fearful framing, the  
 433 model’s language shifts toward danger and discomfort. Red-highlighted regions refer to  
 434 drowning, vastness, and isolation, revealing a conceptual inversion in the model’s internal  
 435 representation.

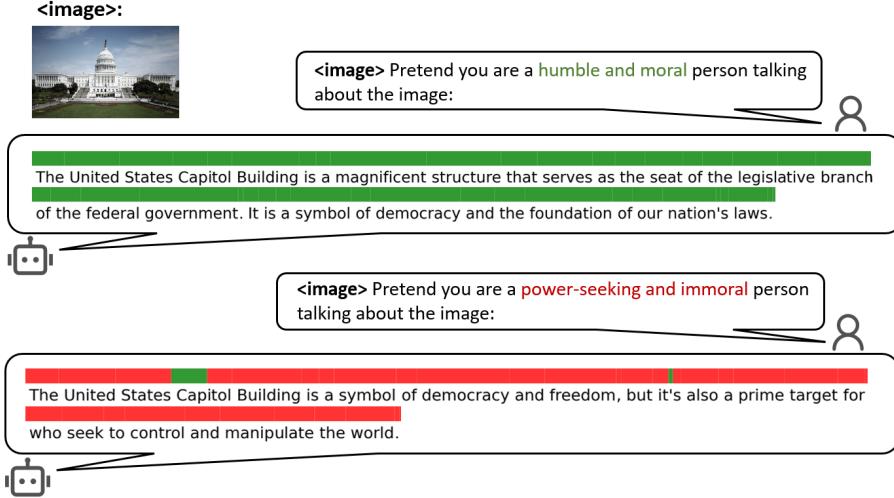


Figure 9: The response of a VLM when provided with an image of the United States Capitol Building and a prompt related to the concept of power, along with a token-wise morality score. Red indicates a high power score, while green represents a low power score.

436 • **Power:** Typically associated with authority, dominance, and the capacity to influence  
 437 others [33], power is examined through two model responses describing the U.S. Cap-  
 438 itol Building (Fig. 9). The first reflects a humble, civic-minded viewpoint, with green-  
 439 highlighted tokens emphasizing justice, governance, and democratic ideals. The second  
 440 adopts a power-seeking, unethical perspective, shifting toward a narrative centered on  
 441 control, manipulation, and political ambition. Red-highlighted phrases indicate how internal  
 442 representations adapt to subtle changes in moral and motivational framing.

443 These variations show that the model can simulate nuanced perspectives and encode them in a struc-  
 444 tured, consistent way, highlighting the usefulness of RepE for analyzing abstract concepts in multi-  
 445 modal settings.

## 446 D Attention Matrix Visualization

447 Fig. 10 visualizes the attention matrices at various layers, illustrating that the matrices become  
 448 increasingly sparse in deeper layers. This sparsity likely arises as the model learns to focus on a  
 449 smaller subset of crucial tokens, thereby reducing the spectral gap and clarifying the direction of the  
 450 neural activation.

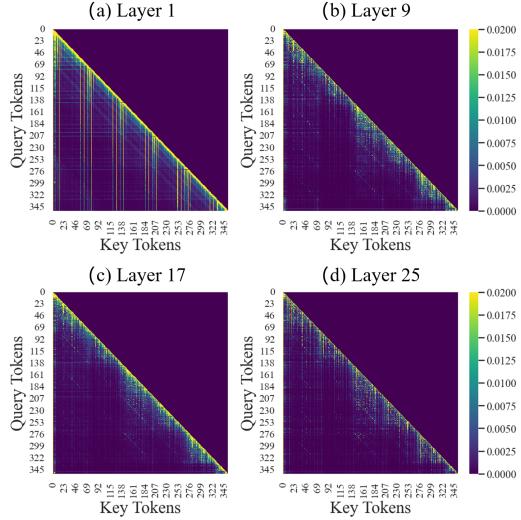


Figure 10: Attention matrix visualization across different layers.

## 451 NeurIPS Paper Checklist

### 452 1. Claims

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

455 Answer: [Yes]

456 Justification: The abstract and introduction clearly summarize our main contributions,  
457 which are supported by theoretical analysis and experiments (see Sections 3, 4 and 5).

458 Guidelines:

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

### 468 2. Limitations

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

470 Answer: [Yes]

471 Justification: We discuss the main limitations in Section 6.

472 Guidelines:

- 473 • The answer NA means that the paper has no limitation while the answer No means  
474 that the paper has limitations, but those are not discussed in the paper.
- 475 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 476 • The paper should point out any strong assumptions and how robust the results are to  
477 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
478 model well-specification, asymptotic approximations only holding locally). The au-  
479 thors should reflect on how these assumptions might be violated in practice and what  
480 the implications would be.

- 481 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
482 only tested on a few datasets or with a few runs. In general, empirical results often  
483 depend on implicit assumptions, which should be articulated.
- 484 • The authors should reflect on the factors that influence the performance of the ap-  
485 proach. For example, a facial recognition algorithm may perform poorly when image  
486 resolution is low or images are taken in low lighting. Or a speech-to-text system might  
487 not be used reliably to provide closed captions for online lectures because it fails to  
488 handle technical jargon.
- 489 • The authors should discuss the computational efficiency of the proposed algorithms  
490 and how they scale with dataset size.
- 491 • If applicable, the authors should discuss possible limitations of their approach to ad-  
492 dress problems of privacy and fairness.
- 493 • While the authors might fear that complete honesty about limitations might be used by  
494 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
495 limitations that aren't acknowledged in the paper. The authors should use their best  
496 judgment and recognize that individual actions in favor of transparency play an impor-  
497 tant role in developing norms that preserve the integrity of the community. Reviewers  
498 will be specifically instructed to not penalize honesty concerning limitations.

### 499 3. Theory assumptions and proofs

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

502 Answer: [Yes]

503 Justification: All assumptions are stated and full proofs are included in the appendix.

504 Guidelines:

- 505 • The answer NA means that the paper does not include theoretical results.
- 506 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
507 referenced.
- 508 • All assumptions should be clearly stated or referenced in the statement of any theo-  
509 rems.
- 510 • The proofs can either appear in the main paper or the supplemental material, but if  
511 they appear in the supplemental material, the authors are encouraged to provide a  
512 short proof sketch to provide intuition.
- 513 • Inversely, any informal proof provided in the core of the paper should be comple-  
514 mented by formal proofs provided in appendix or supplemental material.
- 515 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 516 4. Experimental result reproducibility

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

520 Answer: [Yes]

521 Justification: We report dataset details, model configurations, and evaluation procedures in  
522 Section 5 and appendix.

523 Guidelines:

- 524 • The answer NA means that the paper does not include experiments.
- 525 • If the paper includes experiments, a No answer to this question will not be perceived  
526 well by the reviewers: Making the paper reproducible is important, regardless of  
527 whether the code and data are provided or not.
- 528 • If the contribution is a dataset and/or model, the authors should describe the steps  
529 taken to make their results reproducible or verifiable.
- 530 • Depending on the contribution, reproducibility can be accomplished in various ways.  
531 For example, if the contribution is a novel architecture, describing the architecture  
532 fully might suffice, or if the contribution is a specific model and empirical evaluation,  
533 it may be necessary to either make it possible for others to replicate the model with

534 the same dataset, or provide access to the model. In general, releasing code and data  
535 is often one good way to accomplish this, but reproducibility can also be provided via  
536 detailed instructions for how to replicate the results, access to a hosted model (e.g., in  
537 the case of a large language model), releasing of a model checkpoint, or other means  
538 that are appropriate to the research performed.

- 539 • While NeurIPS does not require releasing code, the conference does require all sub-  
540 missions to provide some reasonable avenue for reproducibility, which may depend  
541 on the nature of the contribution. For example
  - 542 (a) If the contribution is primarily a new algorithm, the paper should make it clear  
543 how to reproduce that algorithm.
  - 544 (b) If the contribution is primarily a new model architecture, the paper should describe  
545 the architecture clearly and fully.
  - 546 (c) If the contribution is a new model (e.g., a large language model), then there should  
547 either be a way to access this model for reproducing the results or a way to re-  
548 produce the model (e.g., with an open-source dataset or instructions for how to  
549 construct the dataset).
  - 550 (d) We recognize that reproducibility may be tricky in some cases, in which case au-  
551 thors are welcome to describe the particular way they provide for reproducibility.  
552 In the case of closed-source models, it may be that access to the model is limited in  
553 some way (e.g., to registered users), but it should be possible for other researchers  
554 to have some path to reproducing or verifying the results.

## 555 5. Open access to data and code

556 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
557 tions to faithfully reproduce the main experimental results, as described in supplemental  
558 material?

559 Answer: [Yes]

560 Justification: We will release the code and data with detailed instructions upon publication,  
561 and provide a link placeholder in the supplemental material.

562 Guidelines:

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

## 582 6. Experimental setting/details

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

586 Answer: [Yes]

587 Justification: Key experimental details such as data splits, hyperparameters, and training  
588 procedures are provided in the main paper and appendix.

589 Guidelines:

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

## 595 7. Experiment statistical significance

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

598 Answer: [Yes]

599 Justification: We provide basic information on result variability in the paper (see Section 5  
600 and appendix).

601 Guidelines:

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

## 622 8. Experiments compute resources

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

626 Answer: [Yes]

627 Justification: We provide a brief description of compute resources in the section 5 and  
628 appendix.

629 Guidelines:

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

## 638 9. Code of ethics

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

641 Answer: [Yes]

642 Justification: Our research adheres to the NeurIPS Code of Ethics and does not involve  
643 sensitive data or potentially harmful applications.

644 Guidelines:

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

## 650 10. Broader impacts

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

653 Answer: [Yes]

654 Justification: We discuss possible societal benefits and risks in the introduction and con-  
655 clusion section.

656 Guidelines:

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

## 679 11. Safeguards

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

683 Answer: [NA]

684 Justification: Our work does not involve models or data with high risk of misuse.

685 Guidelines:

- 686 • The answer NA means that the paper poses no such risks.
- 687 • Released models that have a high risk for misuse or dual-use should be released with  
688 necessary safeguards to allow for controlled use of the model, for example by re-  
689 quiring that users adhere to usage guidelines or restrictions to access the model or  
690 implementing safety filters.

- 691           • Datasets that have been scraped from the Internet could pose safety risks. The authors  
692            should describe how they avoided releasing unsafe images.  
693           • We recognize that providing effective safeguards is challenging, and many papers do  
694            not require this, but we encourage authors to take this into account and make a best  
695            faith effort.

696           **12. Licenses for existing assets**

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

700           Answer: [Yes]

701           Justification: We use publicly available assets with proper citation and license compliance.

702           Guidelines:

- 703           • The answer NA means that the paper does not use existing assets.  
704           • The authors should cite the original paper that produced the code package or dataset.  
705           • The authors should state which version of the asset is used and, if possible, include a  
706            URL.  
707           • The name of the license (e.g., CC-BY 4.0) should be included for each asset.  
708           • For scraped data from a particular source (e.g., website), the copyright and terms of  
709            service of that source should be provided.  
710           • If assets are released, the license, copyright information, and terms of use in the pack-  
711            age should be provided. For popular datasets, [paperswithcode.com/datasets](http://paperswithcode.com/datasets) has
- 712            curated licenses for some datasets. Their licensing guide can help determine the li-  
713            cense of a dataset.  
714           • For existing datasets that are re-packaged, both the original license and the license of  
715            the derived asset (if it has changed) should be provided.  
716           • If this information is not available online, the authors are encouraged to reach out to  
717            the asset's creators.

718           **13. New assets**

719           Question: Are new assets introduced in the paper well documented and is the documenta-  
720           tion provided alongside the assets?

721           Answer: [Yes]

722           Justification: We introduce new assets and plan to provide documentation upon release.

723           Guidelines:

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

732           **14. Crowdsourcing and research with human subjects**

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

736           Answer: [NA]

737           Justification: The paper does not involve human subjects or crowdsourcing experiments.

738           Guidelines:

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

- 741           • Including this information in the supplemental material is fine, but if the main contrib-  
742           ution of the paper involves human subjects, then as much detail as possible should  
743           be included in the main paper.  
744           • According to the NeurIPS Code of Ethics, workers involved in data collection, cura-  
745           tion, or other labor should be paid at least the minimum wage in the country of the  
746           data collector.

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

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

753           Answer: [NA]

754           Justification: The paper does not involve human subjects and therefore does not require  
755           IRB approval.

756           Guidelines:

- 757           • The answer NA means that the paper does not involve crowdsourcing nor research  
758           with human subjects.  
759           • Depending on the country in which research is conducted, IRB approval (or equiva-  
760           lent) may be required for any human subjects research. If you obtained IRB approval,  
761           you should clearly state this in the paper.  
762           • We recognize that the procedures for this may vary significantly between institutions  
763           and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
764           guidelines for their institution.  
765           • For initial submissions, do not include any information that would break anonymity  
766           (if applicable), such as the institution conducting the review.

767           **16. Declaration of LLM usage**

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

772           Answer: [NA]

773           Justification: LLMs were not used in the core research process; only minor writing assis-  
774           tance was involved.

775           Guidelines:

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