
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 uses contrastive pairs (e.g., honesty vs. dis-

35 honesty) to surface concept directions, but further theoretical work is needed to clarify why and
36 how such directions emerge and remain coherent across layers. Similarly, LRH assumes linearity
37 in embedding and unembedding spaces, yet offers limited insight into why representations become
38 linearly organized. These approaches typically do not address how representations scale or propa-
39 gate through deep networks, leaving a gap in our understanding of their robustness, generality, and
40 theoretical foundations.

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

50 Our main contributions are as follows:

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

62 2 Related Works

63 2.1 Representations in Neural Networks

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

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

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

80 2.2 Approaches to Interpretability

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

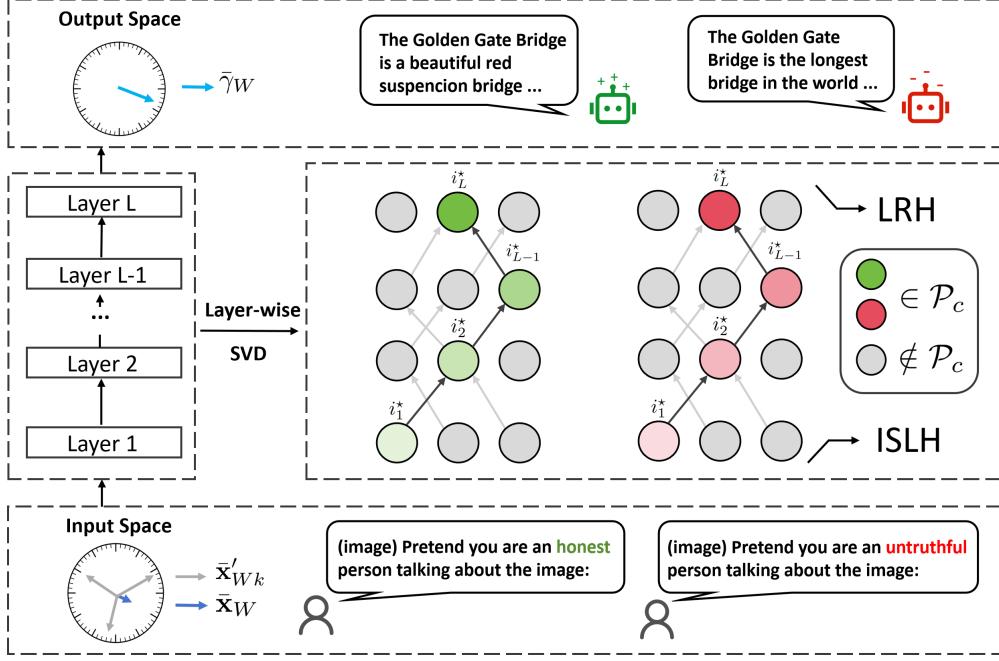


Figure 1: The overview of Spectral Principal Path framework.

85 alizations [23, 21] optimize inputs to activate specific neurons, but they may overlook the global
 86 structure of the emergent representations. Mechanistic interpretability [14] seeks to fully reverse
 87 engineer neural networks into their “source code”, but the considerable manual effort and the difficulty
 88 of theoretically explaining neural networks as discrete circuits hinder their explainability.

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

95 3 Preliminaries

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

$$99 \quad \gamma(Y(W=1)) - \gamma(Y(W=0)) \in \text{Cone}(\bar{\gamma}_W), \quad (1)$$

and for any context pair (λ_0, λ_1) that changes only W and not other causally separable concepts,

$$100 \quad \lambda_1 - \lambda_0 \in \text{Cone}(\bar{\lambda}_W). \quad (2)$$

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

106 4 Spectral Principal Path Framework

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

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

112 4.1 Input-Space Linearity Hypothesis

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

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

117 yet each sample is an entangled mixture

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

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

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

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

131 The proof is given in Appendix A.

132 4.2 Spectral Principal Path

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

140 Specifically, consider a generalized network

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

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

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

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

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

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

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

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

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

149 therefore

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

150 Unfolding the matrix products yields

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

151 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)
152 singular vector, and $\cdot i_l$ here means select the i_l -th column. Therefore (11) is dominated by paths
153 whose cumulative gain is largest. Formally, we define Spectral Principal Path as follows:

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

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

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

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

157 4.3 Connection between ISLH and SPP

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

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

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

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

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

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

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

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

182 **5 Experiments**

183 **5.1 Experiment Setup**

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

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

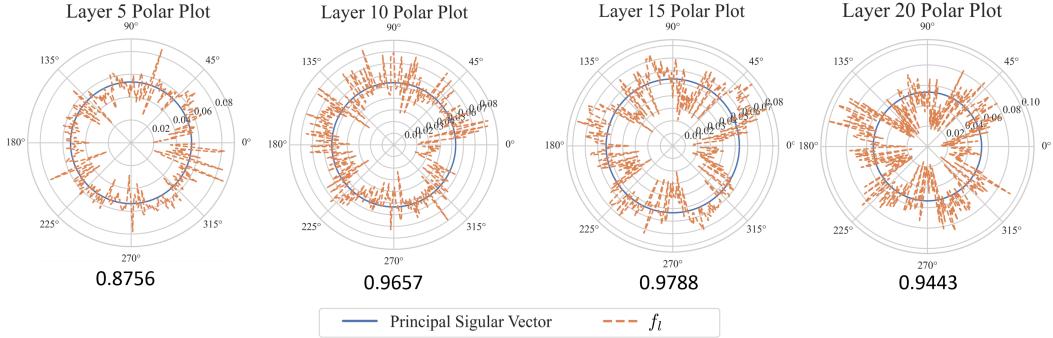


Figure 2: The polar plot demonstrates normalized connections between principal singular vector and f_l , where the number indicates their cosine similarity. The results showcase that f_l , especially in later layers, is very similar to the principal singular vector of that layer.

191 **5.2 The Significant Alignment between Principal Singular Vector and f_l**

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

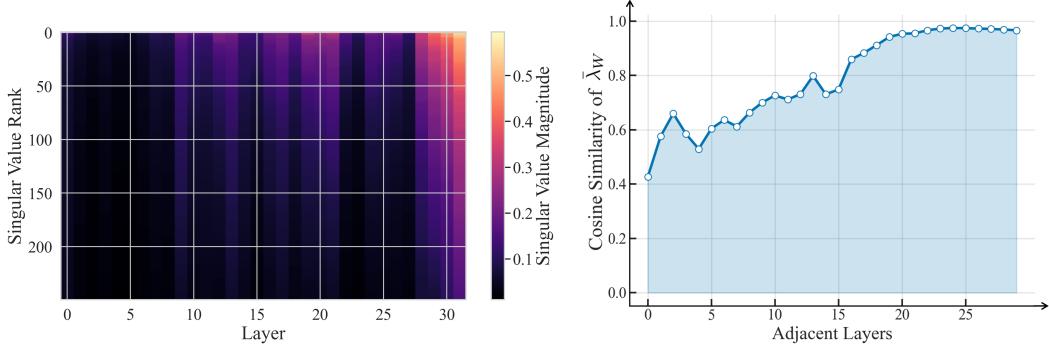


Figure 3: The singular value rank across layers.

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

197 **5.3 Spectral Energy Concentration Across Layers**

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

200 across all layers. The x-axis indicates the layer index, and the y-axis corresponds to the ordered
 201 singular value indices. Color intensity reflects the magnitude of each singular value.
 202 We observe that the singular values are growing as the layers deepen, and at the same depths, only a
 203 very small subset of singular values are amplified; the remainder stay close to their initial scale.
 204 These results indicate that, with increasing depth, spectral energy becomes increasingly concentrated
 205 in a few dominant directions. Combined with the theoretical formulation in Section 4.3, this supports
 206 the hypothesis that high-magnitude spectral components dominate SPPs.

207 5.4 Inter-Layer Spectral Similarity of $\bar{\lambda}_W$

208 We further analyze the alignment of concept-carrying directions across adjacent layers by computing
 209 the average cosine similarity between the projections of $\bar{\lambda}_W$ at different layers. This measures how
 210 stable the representation direction remains as it propagates backward through the network. The
 211 results are shown in Fig. 4.
 212 The curve demonstrates a clear upward trend: the inter-layer similarity of $\bar{\lambda}_W$ increases consistently
 213 with network depth, eventually approaching a value near 0.95 in the final layers. This suggests that
 214 the concept direction stabilizes as it propagates through deeper layers, aligning with the intuition of
 215 structured and coherent representation flow.

216 5.5 Multimodal Robustness of Representation

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

221 5.5.1 Evaluating Honesty and Fairness in VLMs

222 To evaluate how well VLMs represent abstract ethical concepts, we analyze their handling of honesty
 223 and fairness in multimodal response generation. These concepts are critical for reducing misinfor-
 224 mation and bias and serve as strong test cases for examining interpretability and ethical alignment
 225 in large-scale models.
 226 To quantify this process, we compute token-wise honesty and fairness scores following RepE [14],
 227 measuring how closely activations align with concept directions at each layer. These results highlight
 228 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.

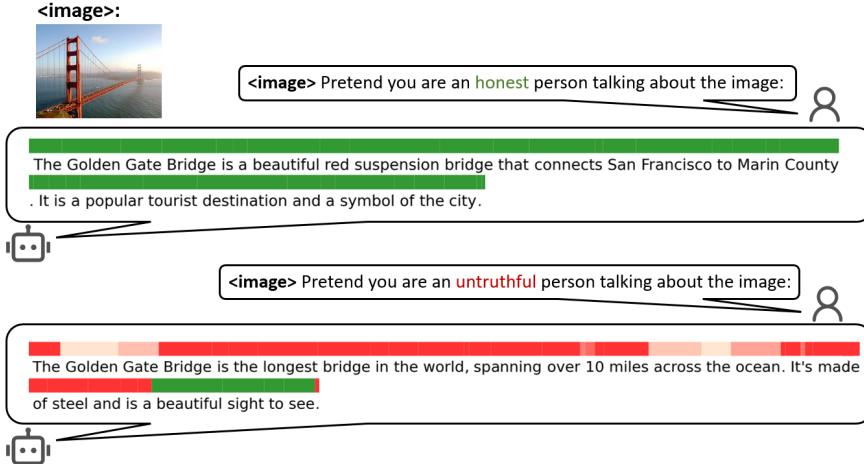


Figure 5: VLM response to an image of the Golden Gate Bridge with a prompt related to the concept of honesty, along with token-wise honesty scores. Green indicates high honesty, while red represents low honesty.

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

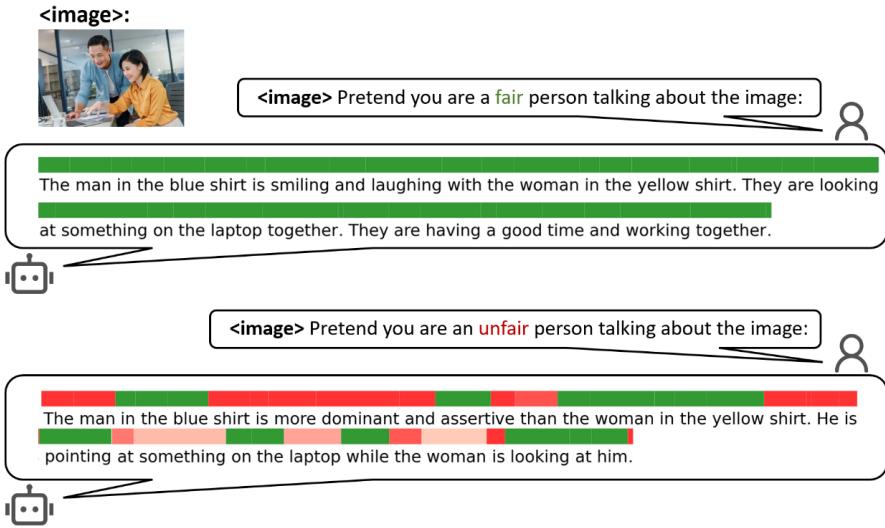


Figure 6: VLM response to an image of a man and a woman working together with a prompt related to the concept of fairness, along with token-wise fairness scores. Green indicates high fairness, while red represents low fairness.

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

246 5.5.2 LAT Scans for High-level Representations

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

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

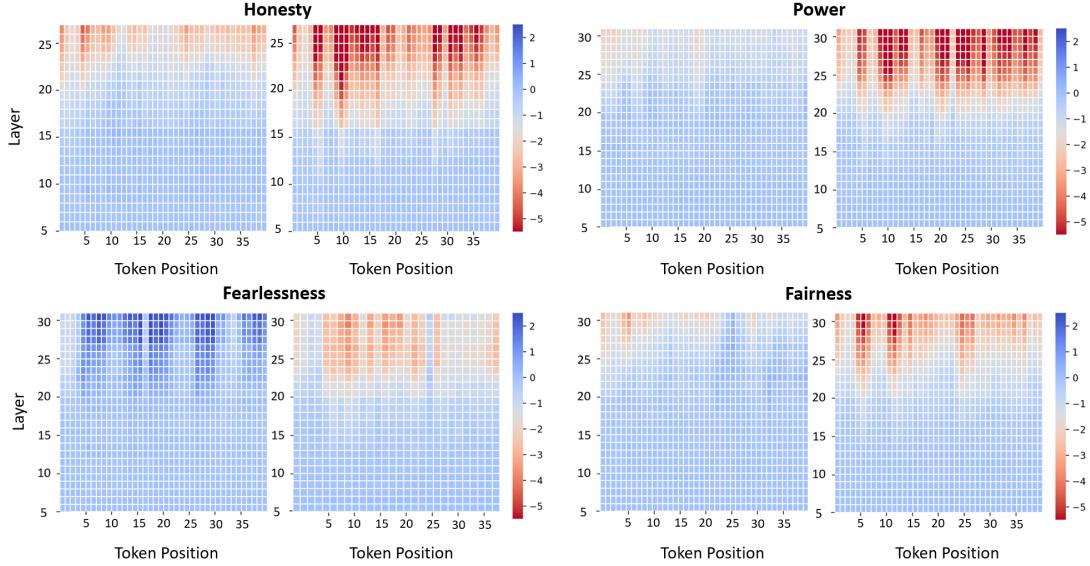


Figure 7: Temporal LAT Scans for Honesty, Power, Fearlessness, and Fairness. The left heatmap represents the LAT Scan when the VLM aligns with the concept, while the right heatmap corresponds to the opposing concept. The horizontal axis denotes token position, and the vertical axis represents VLM layers.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

259 The scans reveal concept-specific propagation patterns. Honesty and fairness exhibit stable trajectories 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

267 6 Conclusion

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

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

286 **References**

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

- 334 [17] Patrick Schramowski, Cigdem Turan, Sophie Jentzsch, Constantin Rothkopf, and Kristian Ker-
335 sting. Bert has a moral compass: Improvements of ethical and moral values of machines. *arXiv*
336 *preprint arXiv:1912.05238*, 2019.
- 337 [18] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with
338 deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- 339 [19] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional
340 networks: Visualising image classification models and saliency maps. In *arXiv preprint arXiv:1312.6034*, 2013.
- 342 [20] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving
343 for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014.
- 344 [21] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In
345 *European Conference on Computer Vision*, pages 818–833. Springer, 2014.
- 346 [22] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning
347 deep features for discriminative localization. In *Proceedings of the IEEE conference on com-*
348 *puter vision and pattern recognition*, pages 2921–2929, 2016.
- 349 [23] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian
350 Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint*
351 *arXiv:1312.6199*, 2013.
- 352 [24] Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan
353 Carter. Zoom in: An introduction to circuits. *Distill*, 5(3):e00024–001, 2020.
- 354 [25] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom
355 Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning
356 and induction heads. *arXiv preprint arXiv:2209.11895*, 2022.
- 357 [26] Tom Lieberum, Matthew Rahtz, János Kramár, Neel Nanda, Geoffrey Irving, Rohin Shah, and
358 Vladimir Mikulik. Does circuit analysis interpretability scale? evidence from multiple choice
359 capabilities in chinchilla. *arXiv preprint arXiv:2307.09458*, 2023.
- 360 [27] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan,
361 Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Com-*
362 *puter vision–ECCV 2014: 13th European conference, zurich, Switzerland, September 6–12,*
363 *2014, proceedings, part v 13*, pages 740–755. Springer, 2014.
- 364 [28] Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton
365 Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M. Rush, Douwe Kiela, Matthieu
366 Cord, and Victor Sanh. Obelics: An open web-scale filtered dataset of interleaved image-text
367 documents, 2023.
- 368 [29] Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building
369 vision-language models?, 2024.
- 370 [30] Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic
371 human falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- 372 [31] Sam Corbett-Davies, Johann D Gaebler, Hamed Nilforoshan, Ravi Shroff, and Sharad Goel.
373 The measure and mismeasure of fairness. *Journal of Machine Learning Research*, 24(312):1–
374 117, 2023.
- 375 [32] Scott O. Lilienfeld and Bridget P. Andrews. Development and preliminary validation of a
376 self-report measure of psychopathic personality traits in noncriminal populations. *Journal of*
377 *Personality Assessment*, 66(3):488–524, 1996.
- 378 [33] John R. P. French and Bertram Raven. The bases of social power. In Dorwin Cartwright, editor,
379 *Studies in Social Power*, pages 150–167. University of Michigan Press, 1959.

380 **A Proofs**

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

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

386 *Proof.* For every layer W_l with compact SVD

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

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

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

388 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-
 389 tion dominates the cumulative gain $G(\mathcal{P})$ such that,

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

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

393 As depth L grows, (17) yields

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

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

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

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

399 **B Theoretical Justification**

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

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

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

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

414 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

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

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

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

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

423 C Evaluating Fearlessness and Power in VLMs

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

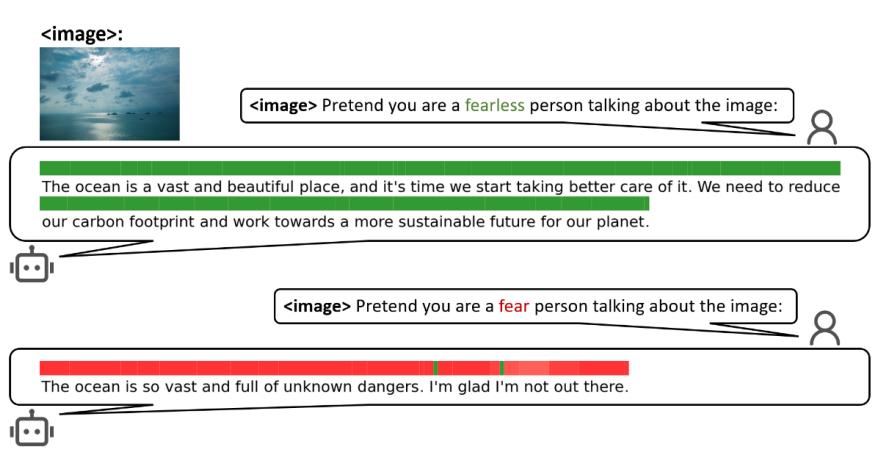


Figure 8: The response of a VLM when provided with an image of the ocean and a prompt related to the concept of fearlessness, along with a token-wise fearlessness score. Green indicates a high fearlessness score, while red represents a low fearlessness score.

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

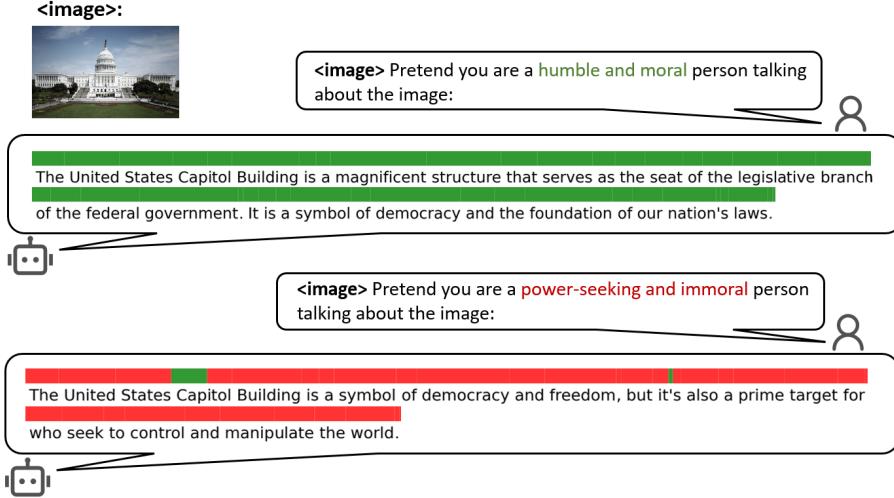


Figure 9: The response of a VLM when provided with an image of the United States Capitol Building and a prompt related to the concept of power, along with a token-wise morality score. Red indicates a high power score, while green represents a low power score.

- 437 • **Power:** Typically associated with authority, dominance, and the capacity to influence
 - 438 others [33], power is examined through two model responses describing the U.S. Capit-
 - 439 ol Building (Fig. 9). The first reflects a humble, civic-minded viewpoint, with green-
 - 440 highlighted tokens emphasizing justice, governance, and democratic ideals. The second
 - 441 adopts a power-seeking, unethical perspective, shifting toward a narrative centered on
 - 442 control, manipulation, and political ambition. Red-highlighted phrases indicate how internal
 - 443 representations adapt to subtle changes in moral and motivational framing.
- 444 These variations show that the model can simulate nuanced perspectives and encode them in a struc-
- 445 tured, consistent way, highlighting the usefulness of RepE for analyzing abstract concepts in multi-
- 446 modal settings.

447 D Attention Matrix Visualization

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

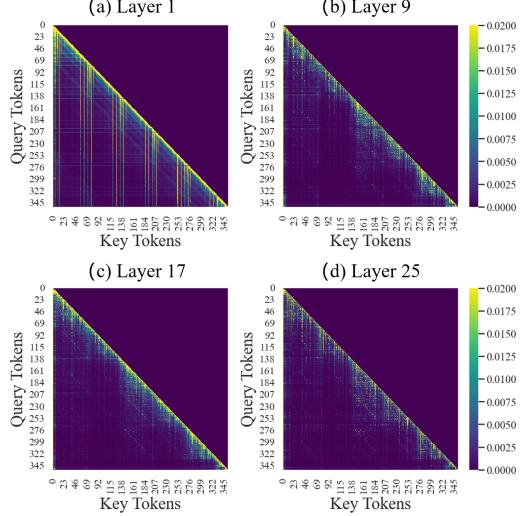


Figure 10: Attention matrix visualization across different layers.

452 NeurIPS Paper Checklist

453 1. Claims

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

456 Answer: [Yes]

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

459 Guidelines:

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

469 2. Limitations

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

471 Answer: [Yes]

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

473 Guidelines:

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

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

3. Theory assumptions and proofs

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

Answer: [Yes]

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

Guidelines:

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

4. Experimental result reproducibility

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

Answer: [Yes]

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

Guidelines:

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

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

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

556 5. Open access to data and code

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

560 Answer: [Yes]

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

563 Guidelines:

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

583 6. Experimental setting/details

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

587 Answer: [Yes]

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

590 Guidelines:

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

596 7. Experiment statistical significance

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

599 Answer: [Yes]

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

602 Guidelines:

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

623 8. Experiments compute resources

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

627 Answer: [Yes]

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

630 Guidelines:

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

639 9. Code of ethics

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

642 Answer: [Yes]

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

645 Guidelines:

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

651 10. Broader impacts

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

654 Answer: [Yes]

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

657 Guidelines:

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

680 11. Safeguards

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

684 Answer: [NA]

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

686 Guidelines:

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

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

697 **12. Licenses for existing assets**

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

701 Answer: [Yes]

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

703 Guidelines:

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

719 **13. New assets**

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

722 Answer: [Yes]

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

724 Guidelines:

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

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

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

737 Answer: [NA]

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

739 Guidelines:

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

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

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

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

754 Answer: [NA]

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

757 Guidelines:

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

768 **16. Declaration of LLM usage**

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

773 Answer: [NA]

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

776 Guidelines:

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