

Mechanistic Interpretability of Vision-Language Models: Tracing Multimodal Concept Binding in CLIP and SigLIP

DR. MADUMERE SMART ONYEMAECHI¹, FRANK UCHEHARA O.²

^{1,2} *Computer Science, Alvan Ikoku Federal University of Education Owerri, Imo State, Nigeria*

Abstract- Vision-Language Models like CLIP and SigLIP align images and text in a shared embedding space, but the internal mechanisms by which they bind visual concepts to language remain opaque. We apply mechanistic interpretability methods to localize and characterize circuits responsible for cross-modal concept binding in CLIP-ViT-B/32 and SigLIP-SO400M. Using causal tracing, sparse autoencoders, and activation patching, we identify three functional subsystems: early visual feature extraction, cross-attention-mediated grounding, and late-stage concept fusion. We find that 12–18% of MLP neurons and 7% of attention heads exhibit polysemantic, modality-invariant concept selectivity for objects, colors, and relations. Causal intervention on these circuits produces predictable changes in image-text alignment scores, confirming causal roles. Our analysis reveals that concept binding relies on a small set of highly interpretable circuits rather than distributed representations. These findings provide a foundation for targeted model editing, bias mitigation, and robust multimodal alignment.

Keywords- mechanistic interpretability, vision-language models, CLIP, concept binding, sparse autoencoders, causal tracing.

I. INTRODUCTION

Vision-Language Models have achieved state-of-the-art performance on zero-shot classification, image captioning, and visual question answering. Models such as CLIP, ALIGN, and SigLIP learn by aligning image and text embeddings through contrastive objectives on web-scale data. Despite their empirical success, these models are black boxes: it is unclear which weights encode concepts like “red car” or “dog chasing ball”, and how visual and linguistic representations are fused.

Mechanistic interpretability seeks to reverse-engineer neural networks into human-understandable algorithms. In language models, this approach has identified circuits for induction heads, factual recall, and arithmetic. Extending these methods to vision-

language models is non-trivial due to modality heterogeneity and the lack of discrete tokens in vision encoders.

We address two questions:

1. Where and how are concepts represented across modalities? Are there neurons or attention heads selective for “cat” in both image and text pathways?
2. Which circuits causally mediate cross-modal binding? Can we intervene to strengthen or break alignment for specific concepts?

Contributions:

1. We adapt causal tracing and activation patching to multimodal settings, enabling causal attribution across vision and text encoders.
2. We train sparse autoencoders on CLIP and SigLIP activations to extract monosemantic features and quantify polysemanticity.
3. We identify and validate circuits for object, attribute, and relational binding, showing that a sparse subset of components controls alignment.
4. We release analysis tools and annotated circuit diagrams for CLIP-ViT-B/32.

II. RELATED WORK

Interpretability in Vision Models: Early work on CNNs used activation maximization and saliency maps. Recent work on ViTs identifies patch-level attention patterns and concept neurons in MLP layers. **Interpretability in Language Models:** Mechanistic studies on GPT-style models uncovered induction heads, knowledge circuits, and task-specific attention patterns. Sparse autoencoders have been used to decompose activations into interpretable features.

Multimodal Interpretability: Prior work on CLIP uses probing, CCA, and concept activation vectors to study alignment. However, these methods are correlational and do not establish causality. Our work

extends causal methods to the multimodal setting and provides circuit-level explanations.

III. METHODS

3.1 Models and Data

We study CLIP-ViT-B/32 and SigLIP-SO400M. For analysis we use 50k image-text pairs from COCO and CC3M, filtered for objects, attributes, and spatial relations. Prompts are of the form “a photo of a {object} with {attribute}”.

3.2 Causal Tracing and Activation Patching

Causal tracing identifies components whose activation states are necessary for correct retrieval. We corrupt the input by zeroing image or text embeddings and restore activations of specific layers or heads. Recovery of the original logit score indicates causal relevance.

Activation patching intervenes by replacing activations from a source input into a target input. We use this to test whether a head or neuron is sufficient to transfer concept binding.

3.3 Sparse Autoencoders for Feature Decomposition

We train k -sparse autoencoders on residual stream activations from MLP and attention outputs. Sparsity $L_0 = 32$, dictionary size 65k. Each SAE feature is scored for polysemanticity and labeled via automated concept probing with 200 ImageNet concepts.

3.4 Evaluation Metrics

- Concept Selectivity: AUROC of a feature’s activation for images/text containing the concept.
- Causal Effect: Change in alignment score after intervention, measured as Δlogit .
- Circuit Sparsity: Fraction of components required to recover 80% of original performance.

IV. RESULTS

4.1 Localization of Concept-Binding Circuits

Causal tracing reveals that early ViT layers encode low-level visual features, while layers 8–12 contain heads that attend to object patches and align them with corresponding text tokens. In the text encoder, layers 6–9 show high causal effect for nouns and adjectives.

We find 1,842 SAE features with high selectivity for objects, 1,203 for colors, and 876 for spatial relations. Only 14% of these features are polysemantic across modalities, suggesting modality-invariant concept encoding.

4.2 Cross-Modal Binding Mechanisms

Two mechanisms dominate:

1. Attention-mediated grounding: Specific heads in layer 10 of the vision encoder attend to object patches and increase alignment when the text contains the matching noun. Patching these heads alone recovers 62% of alignment for “dog” queries.
2. MLP fusion neurons: Neurons in MLP layers 11–12 activate strongly for both image regions and text tokens of the same concept. Ablating these neurons reduces retrieval accuracy by 23% on average.

4.3 Intervention Experiments

Intervening on top 50 concept-selective neurons increases the cosine similarity for the target concept by 0.18 on average, while decreasing similarity for unrelated concepts. Conversely, ablating these neurons causes targeted misalignments without degrading overall performance.

We also observe that relational concepts like “left of” are encoded in smaller, more distributed circuits involving attention heads across both encoders.

V. DISCUSSION

Our results suggest that CLIP-like models implement concept binding using sparse, interpretable circuits rather than fully distributed representations. This aligns with findings in language models but extends them to multimodal fusion.

Implications:

- Model editing: Targeted intervention on concept neurons can debias or update model behavior without full retraining.
- Robustness: Failures in low-resource settings may stem from underdeveloped circuits for specific concepts.
- Evaluation: Causal methods provide a more rigorous evaluation of alignment than probing alone.

Limitations include the focus on CLIP architectures and English text. Extending to decoder-based VLMs and multilingual text is future work.

VI. CONCLUSION

We present a mechanistic analysis of vision-language models, identifying sparse circuits responsible for cross-modal concept binding. Using causal tracing and sparse autoencoders, we show that alignment is mediated by a small set of modality-invariant neurons and attention heads. These results advance understanding of multimodal representation and provide tools for safer, more controllable VLMs.

REFERENCES

- [1] Radford, A., et al. Learning Transferable Visual Models from Natural Language Supervision. ICML, 2021.
- [2] Zhai, X., et al. Sigmoid Loss for Language Image Pre-Training. ICCV, 2023.
- [3] Olah, C., et al. Zoom in: An Introduction to Circuits. Distill, 2020.
- [4] Marks, S., et al. Sparse Feature Circuits: Discovering and Editing Interpretable Causal Graphs in LLMs. arXiv:2403.19647, 2024.
- [5] Cunningham, H., et al. Sparse Autoencoders Find Highly Interpretable Features in Language Models. arXiv:2309.08600, 2023.
- [6] Bau, D., et al. Network Dissection: Quantifying Interpretability of Deep Visual Representations. CVPR, 2017.