Towards Understanding Multimodal Fine-Tuning: A Case Study into Spatial Features

Extended Abstract

- 2 Large vision-language models (VLMs) achieve strong performance by fine-tuning pretrained lan-
- 3 guage backbones to accept projected image tokens alongside text. Yet we lack a mechanistic account
- 4 of how backbone representations adapt and when vision-specific capabilities arise. We introduce a
- 5 stage-wise model diffing framework for multimodal training, extending a method previously applied
- only to language models into the multimodal setting. By comparing sparse autoencoder (SAE) dictio-
- 7 naries across regimes [1], we provide the first mechanistic account of how visual grounding reshapes
- backbone features. Concretely, we warm-start LLaMA-Scope SAEs trained on LLaMA-3.1-8B [2]
- and adapt them to multimodal activations from LLaVA-More using 50k VQAv2 pairs, yielding a
- feature-level view of how a language model develops the ability to "see."
- We finetune SAEs at each transformer block under four regimes: full sequence, image-only, text-only,
- 12 and random initialization, allowing us to separate the effects of vision and text spans. Text-only
- SAEs converge rapidly to near-zero reconstruction error and remain aligned with the base LLM
- dictionary, while image-only and full-sequence runs plateau at higher error due to the mismatch
- between projected image tokens and the pretrained basis. Warm-starting from pretrained SAEs clearly
- outperforms random initialization, stabilizing alignment and showing that finetuning is more effective
- than training from scratch, thereby providing a reliable foundation for stage-wise diffing.
- 18 To detect features reshaped by multimodal fine-tuning, we introduce two signals: (i) modality
- preference, quantified by the variance gap between multimodal and text-only activations, and (ii)
- 20 geometric reorientation, measured by cosine shifts in decoder directions relative to the base SAE.
- 21 A two-stage filter yields a sparse, well-defined set of vision-preferring features that reorient during
- 22 training. This adapted subset represents fewer than 5% of all features, is concentrated in mid-to-deep
- 23 layers, and departs significantly from the original text-only dictionary.
- We then ask whether these adapted features encode spatial grounding. A controlled dataset shift
- 25 from general VQA to spatial queries (e.g., left/right/above/behind) shows that a subset is selectively
- 26 recruited under spatial prompts and remains active under neutral instructions, confirming that they
- 27 capture spatial meaning rather than superficial lexical cues. Automated interpretation and manual
- 28 inspection confirm consistent meanings such as object placement, relative position, and orientation.
- 29 Crucially, we test their causal role. Using attribution patching, a scalable gradient-based proxy for
- 30 activation interventions, we find that a small set of mid to deep attention heads consistently drive
- 31 these spatially selective features. High-scoring heads localize to semantically relevant regions, while
- 32 low-scoring heads fail to capture structure. Ablating either the identified features or the implicated
- 33 heads leads to clear drops in visual spatial reasoning accuracy, showing that these pathways are not
- only correlated with but necessary for spatial grounding.
- 35 Taken together, our findings show that multimodal adaptation is structured and interpretable. A sparse
- 36 set of vision-preferring features reorient to encode spatial relationships, and a correspondingly small
- 37 set of specialized heads act as their causal channel. While our experiments focus on LLaVA-More
- with a LLaMA-3.1-8B backbone, the methodology is general and provides a foundation for auditing
- and refining multimodal training in safety-critical or domain-specific applications.

40 References

- 41 [1] T. Bricken et al. Stage-Wise Model Diffing. Transformer Circuits, 2024. https:// 42 transformer-circuits.pub/2024/model-diffing/index.html
- 43 [2] Z. He et al. *LLaMA-Scope: Extracting millions of features from LLaMA-3.1-8B with sparse autoencoders*.
 44 arXiv:2410.20526, 2024.