
TOWARDS UNDERSTANDING MULTIMODAL FINE-TUNING: A CASE STUDY INTO SPATIAL FEATURES

Anonymous authors

Paper under double-blind review

ABSTRACT

Contemporary Vision–Language Models (VLMs) achieve strong performance on a wide range of tasks by pairing a vision encoder with a pre-trained language model, fine-tuned for visual–text inputs. Yet despite these gains, it remains unclear how language backbone representations adapt during multimodal training and when vision-specific capabilities emerge. In this work, we present the first mechanistic analysis of VLMs adaptation process. Using stage-wise model diffing, a technique that isolates representational changes introduced during multimodal fine-tuning, we reveal how a language model learns to "see". We first identify vision-preferring features that emerge or reorient during fine-tuning. We then show that a selective subset of these features reliably encodes spatial relations, revealed through controlled shifts to spatial prompts. Finally, we trace the causal activation of these features to a small group of attention heads. Our findings show that stage-wise model diffing reveals when and where spatially-grounded multimodal features arise. It also provides a clearer view of modality fusion by showing how visual grounding reshapes features that were previously text-only. This methodology enhances the interpretability of multimodal training and provides a foundation for understanding and refining how pretrained language models acquire vision-grounded capabilities.

1 INTRODUCTION

Large vision–language models (VLMs) have achieved strong performance on multimodal tasks, including visual question answering (VQA), image captioning, object detection, and visual grounding (Li et al., 2024; AI, 2024). These gains are typically realized by fine-tuning pretrained language models to process visual inputs through projected token sequences, allowing for seamless fusion of image and text representations (Xu et al., 2024; Zhang et al., 2024; Dong et al., 2024b;a). Yet we lack a mechanistic account of how language representations adapt during multimodal training and when vision-specific capabilities emerge (Khayatan et al., 2025; Venhoff et al., 2025c; Stan et al., 2024).

In this work, we introduce a method for analyzing multimodal adaptation in VLMs through stage-wise model diffing (Bricken et al., 2024). This mechanistic interpretability technique isolates representational changes introduced during fine-tuning by comparing sparse autoencoder (SAE) dictionaries across training stages, models, or datasets. By tracking how features rotate, emerge, or are repurposed, it has been shown to uncover subtle shifts such as sleeper-agent features (Hubinger et al., 2024). We extend this approach to the multimodal setting, presenting the first application of stage-wise model diffing to study how pretrained language features evolve under visual grounding.

Concretely, we fine-tune LLaMA-Scope SAEs on activations extracted from the LLaVA-More model (He et al., 2024) on 50k VQAv2 dataset samples (Goyal et al., 2017). This warm-start preserves the original feature basis while adapting to multimodal activations. We isolate features that gain visual preference and undergo strong geometric rotation, serving as anchors for studying spatial representations in the backbone. To identify which adapted features encode spatial reasoning, we apply a controlled dataset shift from general VQA to spatial queries. Features that are preferentially recruited under spatial prompts form a selective subset, which we validate through automatic and manual interpretation. These features consistently activate on questions about object placement, relative position, and orientation. Figure 2 highlights the filtered spatial features.

Finally, we use attribution patching to trace the causal pathways by which these spatial features are activated. Our results reveal a sparse set of mid-layer heads that consistently drive spatial

054 representations, often localizing to semantically meaningful regions and reappearing across related
055 prompts. These findings support the hypothesis that a small number of specialized attention heads
056 coordinate visual grounding within the model. Our contributions are as follows:

- 057 • We extend stage-wise model diffing to the multimodal setting, providing the first feature-level
058 account of how pretrained language backbones adapt under visual grounding.
- 059 • We introduce a systematic pipeline to isolate adapted features, identify those selectively
060 recruited by spatial queries, and filter out lexical artifacts.
- 061 • We show that these spatially selective SAE features are functionally involved in reasoning,
062 through empirical evidence and ablation studies, supported by interpretive checks.
- 063 • We causally attribute the emergence of spatial features to a small subset of attention heads
064 using scalable attribution patching, highlighting structured pathways for visual grounding.
065
066

067 By focusing on feature-level change, our approach complements high-level alignment analyses and
068 probing-based methods, providing a deeper mechanistic view of how models “learn to see”. More
069 broadly, this work offers a framework for auditing and refining multimodal training regimes, with
070 implications for safety-critical domains and targeted fine-tuning in specialized applications.
071

072 2 RELATED WORK

073
074 **Model Diffing and Representation Dynamics** Model diffing techniques aim to isolate how internal
075 representations change across models or training stages. Early work focused on coarse similarity
076 measures, such as visualizing function-space geometry (Olah, 2015; Erhan et al., 2010), stitching
077 intermediate layers across models (Lenc & Vedaldi, 2015; Bansal et al., 2021), or defining new
078 similarity metrics (Kornblith et al., 2019; Barannikov et al., 2021). Later studies examined alignment
079 at the level of individual neurons, showing convergent units across independently trained networks
080 (Li et al., 2015; Olah et al., 2020).

081 Sparse autoencoders (SAEs) offered a feature-level lens, and prior work (Kissane et al., 2024) showed
082 that SAEs largely transfer between base and fine-tuned models, implying most features are preserved
083 and only a minority are altered. This motivates methods that can isolate and precisely interpret those
084 changes. Stage-wise model diffing (Bricken et al., 2024) offers such fine-grained resolution, revealing
085 sleeper-agent features and distinguishing between base and chat-tuned models (Minder et al., 2025a).
086 Extensions to multimodal models highlight similar representational shifts, with concept-shift vectors
087 proposed for steering (Khayatan et al., 2025) and evidence that alignment converges in middle-to-late
088 layers (Venhoff et al., 2025c). These remain semantic-level analyses, whereas our work applies
089 stage-wise diffing with SAEs to the backbone, giving the first mechanistic account of multimodal
090 fine-tuning, showing how it rotates features and induces spatial grounding in pretrained language
091 models.

092 **Multimodal Mechanistic Interpretability.** Compared to the rapidly growing literature on mecha-
093 nistic interpretability of textual LLMs, relatively few studies have examined the internal mechanisms
094 of multimodal large language models (MLLMs). Existing work falls into two main categories.

095 First, tool-based and causal analyses aim to explain model behavior at a high level. Approaches
096 include interpretability toolkits based on attention patterns, relevancy maps, and causal interventions
097 (Stan et al., 2024). Other work uses interventions to trace how information is stored and transferred
098 (Basu et al., 2024), or applies causal mediation to study how BLIP integrates visual evidence (Palit
099 et al., 2023). Second, probing-based studies focus on the representations themselves. Several works
100 analyzed CLIP, identifying both strengths and limitations (Tong et al., 2024; Gandelsman et al., 2023;
101 Chen et al., 2023). Others reported multimodal neurons responsive to joint visual–textual concepts
102 (Schwettmann et al., 2023) and examined how VLMs differentiate hallucinated from real objects
103 (Jiang et al., 2024). More recent methods map visual embeddings into linguistic space, projecting
104 features onto language vocabularies (Neo et al., 2024) or showing the late emergence of visual signals
105 in LLM backbones (Venhoff et al., 2025a).

106 In contrast, these studies primarily analyze patterns, interventions, or probing correlations, but do not
107 directly track how multimodal fine-tuning restructures the backbone’s internal features. Our work
addresses this gap by providing a mechanistic perspective.

3 PRELIMINARIES

3.1 VISION-LANGUAGE MODELS

A vision-language model (VLM) consists of a visual encoder f_V , a pretrained language model f_{LM} , and a trainable projector P . The visual encoder (e.g., a ViT (Radford et al., 2021)) extracts image patch embeddings $V = f_V(x) = [v_1, \dots, v_{N_V}]$, which the projector maps into token space $\tilde{V} = P(V)$. These projected image tokens are concatenated with tokenized text embeddings $T = [t_1, \dots, t_{N_T}]$ to form the multimodal sequence $X = [\tilde{v}_1, \dots, \tilde{v}_{N_V}, t_1, \dots, t_{N_T}]$. Alignment between modalities is achieved through *visual instruction tuning*, where image-text pairs fine-tune the backbone to follow multimodal instructions. The language model processes X through transformer layers of multi-head self-attention and feed-forward networks. For each head h , attention is computed as

$$\text{Attn}(Q, K, V) = \text{Softmax}\left(\frac{QK^\top}{\sqrt{d_h}} + M\right)V, \quad (1)$$

where M is the causal mask preventing attention to future tokens. The outputs of all heads are concatenated and projected into the hidden dimension, and mapped through the unembedding matrix to predict next tokens. For our experiments, we adopt LLaVA-More (Cocchi et al., 2025), which extends LLaVA framework (Liu et al., 2023b; 2024) by integrating recent language models and diverse visual backbones; specifically, we use the variant combining the CLIP ViT-Large-Patch14-336 encoder with a LLaMA-3.1-8B language model backbone (Grattafiori et al., 2024).

3.2 SPARSE AUTOENCODERS (SAES)

Sparse Autoencoders (SAEs) learn a dictionary of features that approximate hidden states as sparse linear combinations of interpretable directions, mitigating superposition where many features overlap in the same dimensions (Bricken et al., 2023; Cunningham et al., 2023). Formally, a vanilla SAE encodes $x \in \mathbb{R}^D$ into

$$f(x) = \text{ReLU}(W_{\text{enc}}x + b_{\text{enc}}), \quad \hat{x} = W_{\text{dec}}f(x) + b_{\text{dec}},$$

with $W_{\text{enc}} \in \mathbb{R}^{F \times D}$, $b_{\text{enc}} \in \mathbb{R}^F$, $W_{\text{dec}} \in \mathbb{R}^{D \times F}$, and $b_{\text{dec}} \in \mathbb{R}^D$. Training minimizes

$$\mathcal{L} = \|x - \hat{x}\|_2^2 + \lambda \sum_{i=1}^F |f_i(x)|,$$

combining reconstruction with an L_1 sparsity penalty. Here, decoder columns $(W_{\text{dec}})_{:,i}$ define the direction of each feature in input space, while encoder rows $(W_{\text{enc}})_{i,:}$ act as detectors that determine when a feature is present. Variants such as Top- K SAEs (Gao et al., 2024) further sharpen this tradeoff by enforcing hard sparsity, improving interpretability and reducing feature co-adaptation.

SAEs have been widely applied to uncover monosemantic features and offer a practical lens on model internals, enabling analyses that range from probing knowledge to tracing safety-relevant behaviors (Bricken et al., 2023; Cunningham et al., 2023). They are not, however, a complete decomposition: interpretability can vary across runs and training setups, and recent work suggests their practical utility may be more limited in some settings (Templeton et al., 2024; Kantamneni et al., 2025). Even so, SAEs have proven particularly effective for *model diffing*, where they make it possible to track how features shift across training stages and to surface subtle but behaviorally important dynamics—a direction we expand on in the next subsection (Bricken et al., 2024; Minder et al., 2025b).

3.3 STAGE-WISE MODEL DIFFING

A recent line of work in model diffing has introduced *stage-wise model diffing* (Bricken et al., 2024), which extends SAE analysis across training stages by re-training dictionaries on activations from successive checkpoints while keeping feature indices aligned. This makes it possible to compare whether units are preserved, rotated, or repurposed during adaptation. Applied to controlled fine-tuning trajectories, it disentangles changes due to model updates from dataset shifts and highlights features that drive adaptation. Prior work has shown that stage-wise diffing uncovers fine-grained dynamics, including sleeper-agent features that remain dormant in pretraining but activate once safety constraints are lifted (Hubinger et al., 2024). Compared to crosscoder-based methods (Lindsey et al., 2024), it provides finer resolution at the feature level, though it remains limited to aligned checkpoints of the same architecture.

4 STAGE-WISE MODEL DIFFING FOR MULTIMODAL ADAPTATION

Overview. We aim to understand how multimodal fine-tuning reshapes model representations, using spatial reasoning as a case study of a distinctly multimodal task that integrates both visual and linguistic cues. To this end, we take inspiration from stage-wise diffing 3.3, employing sparse autoencoders (SAEs) as a feature-level lens to track how internal directions shift when a pretrained language backbone is exposed to visual inputs. Our pipeline has three stages. First, we fine-tune SAEs on multimodal activations to obtain a feature dictionary aligned with the vision–language space. Second, we isolate features that prefer visual tokens and undergo substantial geometric rotation, indicating that they have been repurposed by multimodal training. Third, we probe for spatial reasoning by contrasting generic VQA with spatial queries and keeping only features that increase under the shift while remaining active under neutral instructions, ensuring they are not driven by lexical artifacts. In this way, we reduce the original pool of over one million features to a compact set of candidates plausibly recruited for spatial reasoning tasks.

4.1 ADAPTING LANGUAGE DICTIONARIES TO VISION-LANGUAGE SPACE

We start by adapting sparse autoencoders (SAEs) trained on the Llama 3.1 8B backbone to the hidden states of LLAVA-MORE (Llama 3.1 8B backbone) (Cocchi et al., 2025). We use 50k image–question pairs from the VQAv2 dataset (Goyal et al., 2017), a widely used VQA benchmark of images and open-ended questions. Each SAE is attached to the output of a transformer block and trained on cached activations from these samples. Images are represented by 575 consecutive visual tokens, and questions by variable-length text sequences; this separation allows token-type–specific masking.

We initialize SAEs from the pretrained LLAMA-SCOPE release (He et al., 2024), re-instantiated as a Top- K model ($k=50$), preserving a meaningful, interpretable basis. Since our VLM shares the same backbone, this warm-start ensures continuity with the pretrained language feature space and avoids retraining from scratch, allowing us to directly leverage millions of monosemantic features across layers. As a control, we also train SAEs from random initialization under identical conditions. Training uses Adam with a layer-scaled learning rate, and cached activations are processed in padded mini-batches. To disentangle modality-specific contributions, we consider four regimes: (i) full sequence, (ii) image-only, using only the visual-token span, (iii) text-only, using only the non-visual span, and (iv) random initialization. In all cases, the SAE receives the full hidden state sequence, but masking controls which token spans contribute to the training signal.

We evaluate reconstruction quality using the fraction of variance unexplained (FVU) and report sparsity to verify code selectivity. Evaluation is performed on a held-out split. Figure 1 shows FVU as a function of tokens seen across layers and masking regimes. Text-only SAEs converge rapidly, while image-only and full-token regimes converge more slowly to higher error, reflecting the mismatch between projector embeddings and the LLM basis. Random initialization performs worst, underscoring the importance of starting from a pretrained language dictionary. These findings establish text-only SAEs as a reliable reconstruction baseline, which we later use for model diffing.

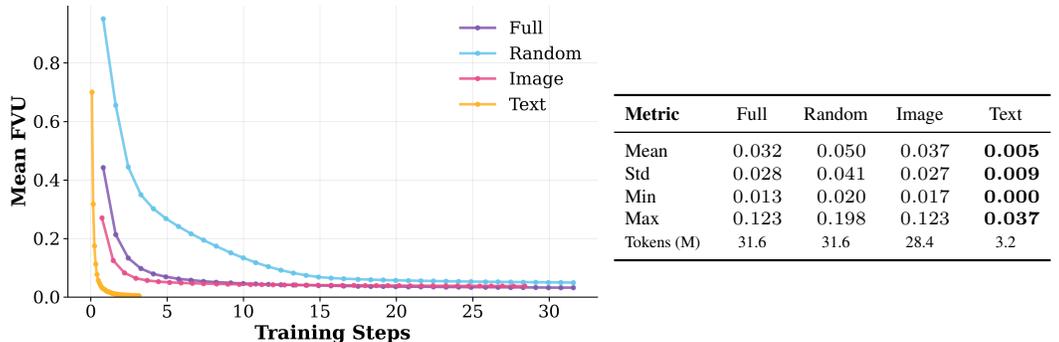


Figure 1: **SAE adaptation on LLAVA-MORE.** Left: Mean fraction of variance unexplained (FVU) across layers on the validation set. Right: Summary statistics of FVU values on the validation set, with decimal alignment; the lowest mean is highlighted in **bold**.

Implications for stage-wise model diffing. Stage-wise diffing assumes that fine-tuning induces *localized* (feature-level) changes rather than wholesale rotations. Prior work reports that image-token representations in early layers exhibit higher reconstruction error than text tokens, indicating a distributional gap between projector outputs and the LLM basis (Venhoff et al., 2025b). Consistent with this, our decoder–cosine analysis (Appx.Fig.5) shows that *text-only* SAEs remain highly aligned to the base LLM dictionary across layers, whereas *image-only* and *full sequence* SAEs undergo large rotations in shallow layers and only align in later layers. We also note that text-only SAEs begin with slightly higher error in the very first layers but adapt extremely quickly, converging to near-zero reconstruction. In contrast, image and full-sequence SAEs plateau at higher error, highlighting the instability of projector-driven spans (see Appx.Fig.6). We therefore focus stage-wise diffing on text-only SAEs, where alignment is stable and feature-level identifiability is more plausible.

4.2 IDENTIFYING ADAPTED FEATURES

We aim to isolate SAE features that (i) undergo geometric reorientation after multimodal adaptation and (ii) show a clear *modality preference* for vision input. Such features are the most informative for model diffing and subsequent causal analysis. To identify them, we rely on two signals:

1. Geometric reorientation (decoder cosine). To test if f has been *repurposed* by multimodal fine-tuning, we compare its decoder direction before and after adaptation. Let $W_{\text{dec},f}^{\text{LLM}}$ be the base SAE decoder vector and $W_{\text{dec},f}^{\text{VLM}}$ the corresponding vector in the VLM-adapted SAE. We compute

$$c_f = \cos(W_{\text{dec},f}^{\text{LLM}}, W_{\text{dec},f}^{\text{VLM}}).$$

High c_f means the semantic direction of f stayed aligned with the original language dictionary; low c_f indicates a substantial rotation, consistent with a reallocation of f to encode new multimodal structure. We use decoder vectors rather than encoder parameters because decoder directions more directly index the feature’s semantics.

2. Modality preference (visual energy). Given the sparsity of SAE activations, we score each feature f by its mean squared activation under vision inputs,

$$E_v(f) = \mathbb{E}_{\text{vision}}[h_f^2],$$

measured on VQA runs of the VLM. Since nearly half of features have $E_v = 0$, a simple cutoff $E_v > \epsilon$ suffices to discard inactive directions and retain those that carry visual signal.

Selection Procedure We define adapted features as those that meet both criteria: $E_v > \epsilon$, ensuring reliable visual responsiveness, and a cosine similarity c_f in the bottom $p_{\text{cos}} = 25\%$, indicating strong decoder rotation. Applying these filters jointly yields a globally defined set comprising about 5% of all features. The joint distribution of E_v and c_f is shown in Fig. 2, with the selected subset highlighted in pink. Details on threshold choices, together with per-layer counts and mean cosine similarities, are provided in Appx. Fig. 7a and Appx. 7b.

4.3 CASE STUDY: IDENTIFYING SPATIAL REASONING FEATURES

We identify spatial features using two signals: (i) recruitment under a shift to spatial queries, and (ii) persistence under neutral prompts that rule out lexical artifacts.

Datasets. Our analysis uses two evaluation sets from VQAv2. The baseline is the full validation split, denoted $\mathcal{D}_{\text{base}}$. To induce a targeted shift, we construct a spatial subset \mathcal{D}_{sp} by filtering questions that contain spatial cues (e.g., *left/right/above/behind*). This contrast tests whether some SAE features are selectively recruited under spatial reasoning.

1. Distribution shift Let $h_f(x_t) \geq 0$ denote the activation of feature f on token t of input x . For a dataset \mathcal{D} , the firing frequency of f is

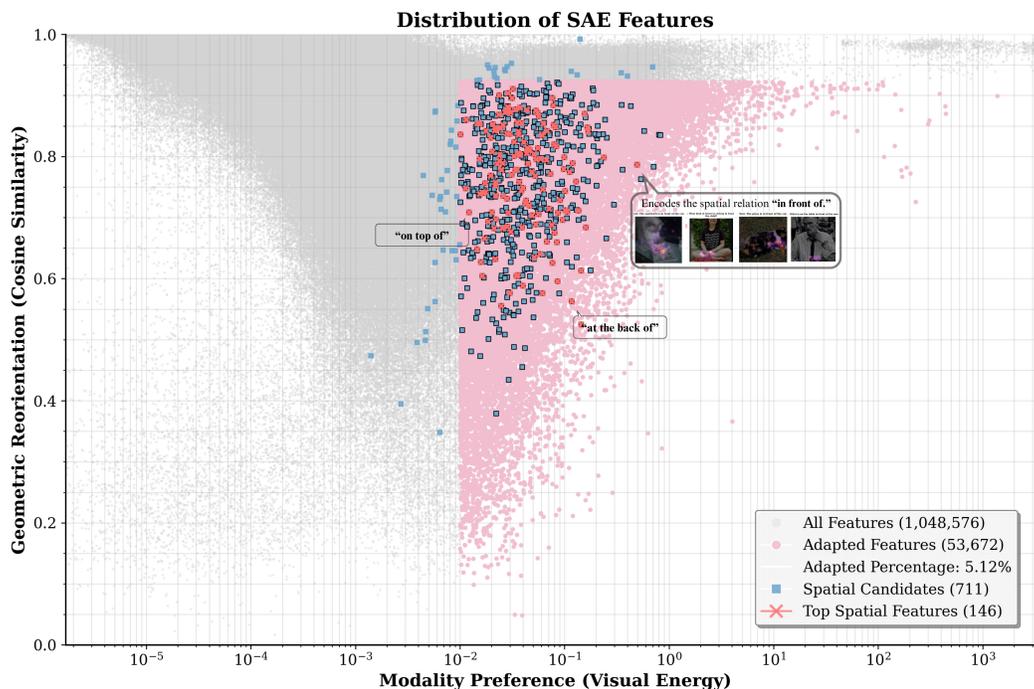
$$p_f(\mathcal{D}) = \frac{1}{n(\mathcal{D})} \sum_{x \in \mathcal{D}} \sum_t \mathbf{1}\{h_f(x_t) > 0\},$$

where $n(\mathcal{D})$ is the total number of tokens. We compute this measure for the base split $\mathcal{D}_{\text{base}}$ and a spatial split \mathcal{D}_{sp} , and evaluate each feature using the frequency gap $\Delta p_f = p_f(\mathcal{D}_{\text{sp}}) - p_f(\mathcal{D}_{\text{base}})$

270 alongside its odds ratio OR_f . Features with meaningful Δp_f and OR_f are flagged as spatial
 271 *candidates* in Fig.2. Further details, including firing-frequency and scatter-plot visualizations for
 272 both splits, are provided in Appx. A.5.
 273
 274

275 **2. Filtering lexical artifacts.** To rule out prompt-lexical effects, we replace the original questions in
 276 each top-activating sample with neutral spatial prompts such as “Describe the positions of objects in
 277 the image.”. Features that continue firing under these generic instructions are preserved as genuinely
 278 image-grounded, while those that fail to activate are discarded. This ensures that the surviving units
 279 reflect spatial reasoning rather than memorized lexical cues.

280 From these filtered candidates, we retain only those also in the adapted set \mathcal{A} (Sec. 4.2), ensuring they
 281 reorient under multimodal fine-tuning and respond to spatial shifts. The surviving features are shown
 282 in Fig. 2 (blue). A subset, marked with red crosses, is further analyzed via automated interpretation,
 283 attribution patching, and ablations (Sec. 5.1, 5.2).
 284



308 Figure 2: **Distribution of SAE features by visual energy and cosine similarity.** All features are
 309 shown in gray; adapted features are highlighted in pink. Spatial candidates are marked with blue
 310 squares, and the subset used for downstream analysis is shown as red crosses.
 311

312

313

314 **Extension to OCR-style prompts.** While our primary case study focuses on spatial reasoning, the
 315 same feature-selection procedure can be applied to other visually grounded skills. As a second case
 316 study, we analyze features associated with visual text recognition by contrasting OCR-style prompts
 317 (e.g., “What does the sign say?”) with generic VQA questions. We construct an OCR-focused split by
 318 filtering VQAv2 images that contain legible embedded text and computing feature firing frequencies
 319 under the same distribution-shift statistics used for spatial queries. This reveals a compact subset
 320 of adapted units whose activations increase on OCR prompts and remain non-zero under neutral
 321 image descriptions, indicating that they are tied to image-grounded text rather than specific lexical
 322 patterns. Additional qualitative examples and follow-up analyses are provided in Appx. A.6, where
 323 these OCR-selective features are shown to align with regions containing characters and words and to
 be supported by a small number of recurring mid-layer heads.

5 EXPERIMENTS

5.1 AUTO-INTERP AND PRELIMINARY INSPECTION

As an initial step toward understanding the selected features, we carried out a preliminary inspection using an automated interpretation pipeline. For each feature, we collect its top-activating samples from two sources: general VQA questions from VQAv2 (not restricted to spatial reasoning) and the Visual Spatial Reasoning (VSR) dataset (Liu et al., 2023a), which is inherently spatial. This pairing allows us to check whether the same underlying meaning emerges consistently across both settings (Fig. 3). A subset of the combined samples are then passed to the `gpt-4o-mini` (OpenAI, 2024) API, which proposes a concise one-sentence description for each feature and assigns an interpretability confidence score based on F1 from a validation classification task. The resulting outputs are stored together with the selection metrics from Sec. 4.2, and are lightly reviewed by hand, so that the retained set reflects both automatic labeling and human verification (see App. B.1 for additional examples and scoring details).

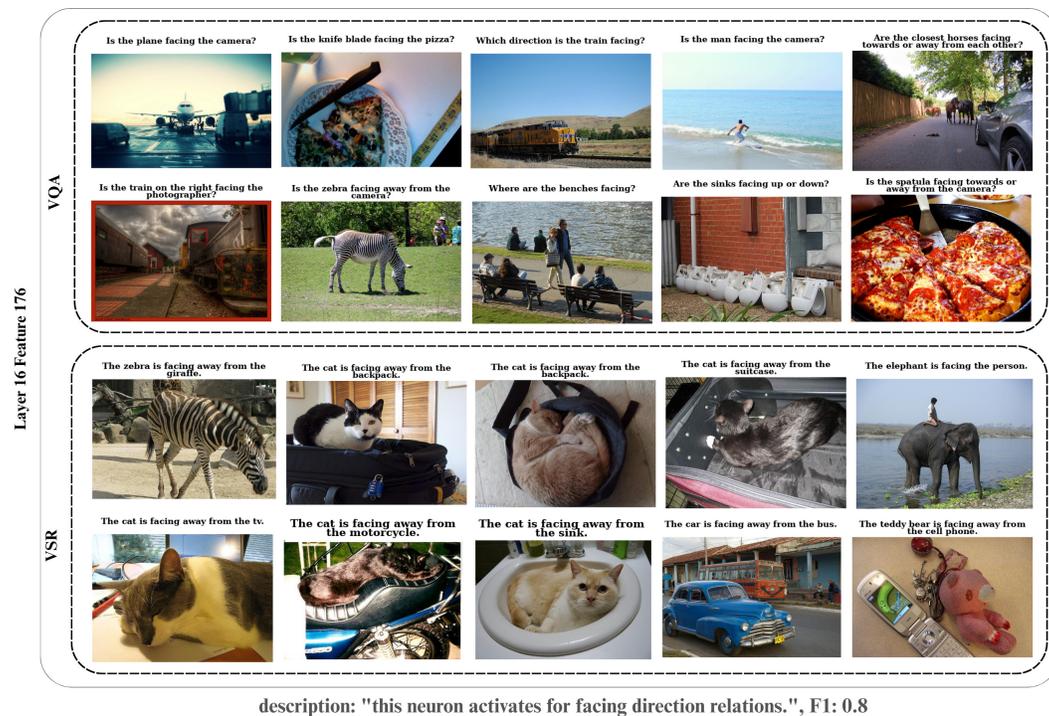


Figure 3: **Auto-Interp example (Layer 16, Feature 176)**. Top VQA and VSR samples both highlight *facing direction*, with activation on objects described as facing toward, away, or relative to another.

5.2 ATTRIBUTION PATCHING TO IDENTIFY SPATIAL HEADS

Method. Attribution patching (Nanda, 2023) is a scalable alternative to activation patching (Zhang & Nanda, 2024), which measures causal effects by replacing activations with counterfactuals. While activation patching requires a separate forward pass per intervention, attribution patching uses a gradient-based linear approximation to estimate interventions with two forward and one backward pass. This makes it practical to probe attribution scores across layers and heads in MLLMs.

We adapt attribution patching to identify which attention heads drive spatially selective SAE features. For a target feature f at layer L , we define a scalar objective by projecting the layer- L activations onto the SAE decoder vector. Gradients of this objective w.r.t. upstream query/key activations indicate how strongly each attention head contributes to f . We compare two runs:

- **Clean run:** the original image–text input.

- **Corrupt run:** the same input, but with layer-0 visual token embeddings replaced by a *mean embedding* computed over many VQA samples. This corruption preserves plausible distributional statistics while deliberately suppressing spatial information.

We then compute two attribution variants, differing in whether the perturbation direction is taken from the corrupted or the clean representation:

$$\text{Method A: } (\text{corr} - \text{clean}) \cdot \nabla_{\text{clean}},$$

$$\text{Method B: } (\text{clean} - \text{corr}) \cdot \nabla_{\text{corr}}.$$

Method A measures how strongly the clean gradients indicate that ablating spatial detail affects the feature, whereas Method B measures how strongly the corrupted gradients indicate that retaining spatial detail matters. In both cases, we obtain per-layer and per-head attribution scores, averaged over the top- k VQA samples that most strongly activate f .

Results. Across the spatially selective features we examined, attribution patching with both methods reveals consistent trends. Layer-wise attribution curves typically peak in middle layers, consistent with the emergence of spatial features in Sec. 4.3 (Appx. Fig. 12). At the head level, both methods generally highlight a small subset of heads with notably high scores, and the top heads identified are often consistent across the two attribution methods (Appx. Fig. 13). This suggests that spatial information is mediated by a specialized group of heads rather than being spread uniformly across the model.

To illustrate the effect of attribution patching on individual features, Appx. Fig. 14 provides detailed examples. In each case, attribution scores isolate a handful of heads, and qualitative maps confirm that high-scoring heads focus on regions consistent with the queried relation (e.g., “on top of,” “behind”), whereas low-scoring heads fail to do so. These head-level overlays can also be used to (i) improve the confidence of automated feature interpretation by coupling sample activations with attention visualizations, and (ii) examine failure cases by checking whether the top spatial features and heads attend to valid regions in misclassified samples. Interestingly, when we look across multiple related spatial features together, we find that some of the same heads recur across related spatial relations. Fig 4 illustrates this pattern. In the top row, L13H1 attends to semantically relevant regions across queries. As a control, the middle row shows that bottom-ranked heads on the same samples fail to localize meaningfully. The bottom row further confirms that irrelevant queries do not trigger spurious activation. More generally, these same heads also attend to meaningful regions such as salient objects or attributes under custom prompts (Appx. Fig. 16), underscoring that attribution patching identifies a set of heads that reliably carry spatial-semantic signal.

5.3 ABLATION STUDY

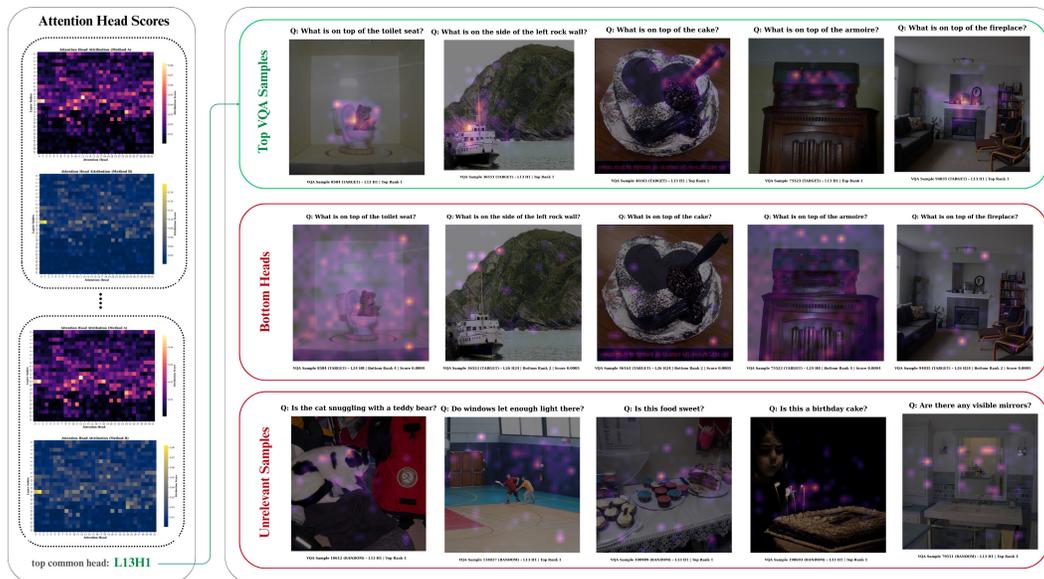
We test whether adapted SAE features are *causally involved* in spatial reasoning by ablating them during inference and measuring performance on VSR (Liu et al., 2023a), a dataset of text-image pairs spanning dozens of spatial relations, and on a *Yes/No* subset of VQA_{v2} (general). Each feature is evaluated on a *relation-specific subset* of VSR constructed from its top-activating samples, so that the ablation directly targets the relation it most strongly encodes. To ablate a target feature f at layer L , we orthogonally remove its decoder direction v (unit norm) from the residual stream at *text* token positions, leaving image tokens unchanged:

$$y \leftarrow y - (y^\top v) v.$$

Evaluation metrics. We report: (i) accuracy drop on VSR (Δ VSR Acc; \downarrow is worse), (ii) accuracy drop on VQA (Δ VQA Acc), (iii) accuracy drop from ablating same-layer random features (Δ Ctrl), and (iv) odds ratio under the spatial distribution shift (VSR OR; \uparrow is better). All runs use identical cached indices, and results are averaged over seeds.

Interpretation. Ablating the top spatial features lowers VSR accuracy by 9–16 points on average while leaving general VQA nearly unchanged (≤ 1 pp), indicating that these directions are functionally used for spatial reasoning rather than general behavior. This shows that probing or switching off a single feature can selectively disable spatial reasoning without harming overall ability. High odds ratios further show selective recruitment under spatial prompts. Random-feature controls yield effects near zero or inconsistent in sign, supporting specificity. Full per-feature results, probability deltas, and seed-wise summaries are reported in Appx. D.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450



451 Figure 4: Attribution patching across related spatial features. Top: recurring top-scoring head
452 (L13H1) localizes to relevant regions in queries about “on top of” relations. Middle: bottom-ranked
453 heads on the same samples fail to capture spatial structure. Bottom: unrelated queries confirm that
454 the top head does not spuriously activate.

455
456
457
458
459
460
461
462

Layer	Feature	Δ VSR Acc (\downarrow)	Δ VQA Acc (\downarrow)	Δ Ctrl	VSR Relations	VSR OR (\uparrow)
7	15870	-15.54	-0.10	-0.88	above	4.32
11	27061	-12.77	-0.40	0.00	across from	8.03
9	15404	-11.19	-0.80	1.08	below	5.60
14	17873	-10.21	-0.30	-1.71	at the right side of	7.17
12	23874	-9.05	-0.40	-0.95	left of	9.10
18	29948	-7.98	-0.30	0.00	beside	8.36

463
464
465
466
467

Table 1: **Top ablated SAE features** ranked by VSR accuracy drop. Columns 2–4 show accuracy
464 drops on VSR, VQA, and random-feature controls; the final column gives odds ratio (VSR OR)
465 as a measure of selective recruitment. Large Δ VSR Acc with small Δ VQA Acc indicates spatial
466 specificity, while near-zero Δ Ctrl confirms robustness.

468
469
470

6 LIMITATIONS

471
472
473
474
475

Our analyses indicate spatial selectivity, but more detailed ablation and steering studies are needed to
472 fully validate causality. Moreover, our experiments are limited to a single model (LLaVA-More with
473 a LLaMA-3.1-8B backbone); applying the method to other backbones and larger corpora will be key
474 to assessing generality.

476
477
478

7 CONCLUSION

479
480
481
482
483
484
485

We set out to understand how a pretrained language backbone learns to “see” under multimodal
480 fine-tuning. By extending stage-wise model diffing to the vision–language setting, we isolated
481 vision-preferring features that undergo strong rotations during training, showed that a subset reliably
482 encodes spatial relations, and traced their causal drivers to a small number of mid-layer attention
483 heads. These results show that multimodal adaptation is structured and interpretable as it can be
484 localized, probed, and explained at the feature level. Beyond spatial reasoning, our methodology
485 offers a general framework for uncovering when and where new capabilities emerge in large models,
showing that multimodal adaptation follows structured patterns rather than diffuse changes. We view

486 this work as an early step toward a mechanistic science of multimodal training, where models can be
487 interpreted both in terms of their outputs and the internal features that support them.
488

489 REFERENCES

- 491 Mistral AI. Pixtral 12b: A new frontier in image and text understanding. [https://mistral.](https://mistral.ai/news/pixtral-12b/)
492 [ai/news/pixtral-12b/](https://mistral.ai/news/pixtral-12b/), September 2024. Accessed: 2024-12-21.
- 493 Yamini Bansal, Preetum Nakkiran, and Boaz Barak. Revisiting model stitching to compare neural
494 representations. *Advances in neural information processing systems*, 34:225–236, 2021.
495
- 496 Serguei Barannikov, Ilya Trofimov, Nikita Balabin, and Evgeny Burnaev. Representation topol-
497 ogy divergence: A method for comparing neural network representations. *arXiv preprint*
498 *arXiv:2201.00058*, 2021.
- 499 Samyadeep Basu, Martin Grayson, Cecily Morrison, Besmira Nushi, Soheil Feizi, and Daniela
500 Massiceti. Understanding information storage and transfer in multi-modal large language models.
501 *Advances in Neural Information Processing Systems*, 37:7400–7426, 2024.
502
- 503 Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick
504 Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec,
505 Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina
506 Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and
507 Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary
508 learning. *Transformer Circuits Thread*, 2023. [https://transformer-circuits.pub/2023/monosemantic-](https://transformer-circuits.pub/2023/monosemantic-features/index.html)
509 [features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html).
- 510 Trenton Bricken, Siddharth Mishra-Sharma, Jonathan Marcus, Adam Jermyn, Christopher Olah,
511 Kelley Rivoire, and Thomas Henighan. Stage-wise model diffing. 2024. [https://](https://transformer-circuits.pub/2024/model-diffing/index.html)
512 transformer-circuits.pub/2024/model-diffing/index.html.
- 513 Haozhe Chen, Junfeng Yang, Carl Vondrick, and Chengzhi Mao. Interpreting and controlling vision
514 foundation models via text explanations. *arXiv preprint arXiv:2310.10591*, 2023.
- 515 Federico Cocchi, Nicholas Moratelli, Davide Caffagni, Sara Sarto, Lorenzo Baraldi, Marcella Cornia,
516 and Rita Cucchiara. Llava-more: A comparative study of llms and visual backbones for enhanced
517 visual instruction tuning. *arXiv preprint arXiv:2503.15621*, 2025.
518
- 519 Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoen-
520 coders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*,
521 2023.
- 522 Guanting Dong, Chenghao Zhang, Mengjie Deng, Yutao Zhu, Zhicheng Dou, and Ji-Rong
523 Wen. Progressive multimodal reasoning via active retrieval. 2024a. URL [https://api.](https://api.semanticscholar.org/CorpusID:274859457)
524 [semanticscholar.org/CorpusID:274859457](https://api.semanticscholar.org/CorpusID:274859457).
- 525 Yuhao Dong, Zuyan Liu, Hai-Long Sun, Jingkang Yang, Winston Hu, Yongming Rao, and Ziwei Liu.
526 Insight-v: Exploring long-chain visual reasoning with multimodal large language models. *arXiv*
527 *preprint arXiv:2411.14432*, 2024b.
528
- 529 Dumitru Erhan, Yoshua Bengio, Aaron Courville, Pierre-Antoine Manzagol, Pascal Vincent, and
530 Samy Bengio. Why does unsupervised pre-training help deep learning? 11:625–660, March 2010.
531 ISSN 1532-4435.
- 532 Yossi Gandelsman, Alexei A Efros, and Jacob Steinhardt. Interpreting clip’s image representation via
533 text-based decomposition. *arXiv preprint arXiv:2310.05916*, 2023.
534
- 535 Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya
536 Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. *arXiv preprint*
537 *arXiv:2406.04093*, 2024.
- 538 Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa
539 matter: Elevating the role of image understanding in visual question answering. In *Proceedings of*
the IEEE conference on computer vision and pattern recognition, pp. 6904–6913, 2017.

540 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad
541 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of
542 models. *arXiv preprint arXiv:2407.21783*, 2024.

543
544 Zhengfu He, Wentao Shu, Xuyang Ge, Lingjie Chen, Junxuan Wang, Yunhua Zhou, Frances Liu,
545 Qipeng Guo, Xuanjing Huang, Zuxuan Wu, et al. Llama scope: Extracting millions of features
546 from llama-3.1-8b with sparse autoencoders. *arXiv preprint arXiv:2410.20526*, 2024.

547 Evan Hubinger, Carson Denison, Jesse Mu, Mike Lambert, Meg Tong, Monte MacDiarmid, Tamera
548 Lanham, Daniel M Ziegler, Tim Maxwell, Newton Cheng, et al. Sleeper agents: Training deceptive
549 llms that persist through safety training. *arXiv preprint arXiv:2401.05566*, 2024.

550 Nick Jiang, Anish Kachinthaya, Suzie Petryk, and Yossi Gandelsman. Interpreting and editing
551 vision-language representations to mitigate hallucinations. *arXiv preprint arXiv:2410.02762*, 2024.

552
553 Subhash Kantamneni, Joshua Engels, Senthooan Rajamanoharan, Max Tegmark, and Neel Nanda.
554 Are sparse autoencoders useful? a case study in sparse probing. *arXiv preprint arXiv:2502.16681*,
555 2025.

556 Pegah Khayatan, Mustafa Shukor, Jayneel Parekh, and Matthieu Cord. Analyzing fine-tuning
557 representation shift for multimodal llms steering alignment. *arXiv preprint arXiv:2501.03012*,
558 2025.

559
560 Connor Kissane, Robert Krzyzanowski, Arthur Conmy, and Neel Nanda. Saes (usually) transfer
561 between base and chat models. AI Alignment Forum post, July 18 2024.

562 Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural
563 network representations revisited. In *International conference on machine learning*, pp. 3519–3529.
564 PMIR, 2019.

565
566 Karel Lenc and Andrea Vedaldi. Understanding image representations by measuring their equivariance
567 and equivalence. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
568 pp. 991–999, 2015.

569 Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li.
570 Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv*
571 *preprint arXiv:2407.07895*, 2024.

572
573 Yixuan Li, Jason Yosinski, Jeff Clune, Hod Lipson, and John Hopcroft. Convergent learning: Do
574 different neural networks learn the same representations? *arXiv preprint arXiv:1511.07543*, 2015.

575 Jack Lindsey, Anh Tuan Tran, Neel Nanda, Trenton Bricken, Adam Jermyn, Kelley Rivoire, and
576 Chris Olah. Sparse crosscoders for cross-layer features and model diffing. *Transformer Circuits*
577 *Thread*, 2024. URL [https://transformer-circuits.pub/2024/crosscoders/](https://transformer-circuits.pub/2024/crosscoders/index.html)
578 [index.html](https://transformer-circuits.pub/2024/crosscoders/index.html).

579 Fangyu Liu, Guy Edward Toh Emerson, and Nigel Collier. Visual spatial reasoning. *Transactions of*
580 *the Association for Computational Linguistics*, 2023a.

581
582 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in*
583 *neural information processing systems*, 36:34892–34916, 2023b.

584
585 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction
586 tuning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
587 pp. 26296–26306, 2024.

588
589 Julian Minder, Clément Dumas, Caden Juang, Bilal Chughtai, and Neel Nanda. Robustly identifying
590 concepts introduced during chat fine-tuning using crosscoders. *arXiv preprint arXiv:2504.02922*,
591 2025a.

592
593 Julian Minder, Clément Dumas, and Neel Nanda. What we learned try-
ing to diff base and chat models (and why it matters). *LessWrong*, 2025b.
URL [https://www.lesswrong.com/posts/xmpauEXEerzYcJKNm/](https://www.lesswrong.com/posts/xmpauEXEerzYcJKNm/what-we-learned-trying-to-diff-base-and-chat-models-and-why)
[what-we-learned-trying-to-diff-base-and-chat-models-and-why](https://www.lesswrong.com/posts/xmpauEXEerzYcJKNm/what-we-learned-trying-to-diff-base-and-chat-models-and-why).

594 Neel Nanda. Attribution patching: Activation patching at industrial scale. [https://www.](https://www.neelnanda.io/mechanistic-interpretability)
595 [neelnanda.io/mechanistic-interpretability](https://www.neelnanda.io/mechanistic-interpretability), 2023. Accessed: 2025-08-23.

596

597 Clement Neo, Luke Ong, Philip Torr, Mor Geva, David Krueger, and Fazl Barez. Towards interpreting
598 visual information processing in vision-language models. *arXiv preprint arXiv:2410.07149*, 2024.

599

600 Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter.
601 Zoom in: An introduction to circuits. *Distill*, 2020. doi: 10.23915/distill.00024.001.
602 <https://distill.pub/2020/circuits/zoom-in>.

603

604 Christopher Olah. Visualizing representations: Deep learning and human beings. [https://colah.](https://colah.github.io/posts/2015-01-Visualizing-Representations/)
605 [github.io/posts/2015-01-Visualizing-Representations/](https://colah.github.io/posts/2015-01-Visualizing-Representations/), 2015. Accessed:
606 2025-08-23.

607

608 OpenAI. Gpt-4o-mini: Advancing cost-efficient intelligence. [https://openai.com/index/](https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/)
609 [gpt-4o-mini-advancing-cost-efficient-intelligence/](https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/), 2024. Accessed:
610 2024-12-21.

611

612 Vedant Palit, Rohan Pandey, Aryaman Arora, and Paul Pu Liang. Towards vision-language mechanistic
613 interpretability: A causal tracing tool for blip. In *Proceedings of the IEEE/CVF International*
614 *Conference on Computer Vision*, pp. 2856–2861, 2023.

615

616 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
617 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
618 models from natural language supervision. In *International conference on machine learning*, pp.
619 8748–8763. PmLR, 2021.

620

621 Sarah Schwettmann, Neil Chowdhury, Samuel Klein, David Bau, and Antonio Torralba. Multimodal
622 neurons in pretrained text-only transformers. In *Proceedings of the IEEE/CVF International*
623 *Conference on Computer Vision*, pp. 2862–2867, 2023.

624

625 Gabriela Ben Melech Stan, Estelle Aflalo, Raanan Yehezkel Rohekar, Anahita Bhiwandiwalla,
626 Shao-Yen Tseng, Matthew Lyle Olson, Yaniv Gurwicz, Chenfei Wu, Nan Duan, and Vasudev
627 Lal. Lvlm-interpret: an interpretability tool for large vision-language models. *arXiv preprint*
628 *arXiv:2404.03118*, 2024.

629

630 Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen,
631 Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L
632 Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers,
633 Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan.
634 Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Trans-*
635 *former Circuits Thread*, 2024. URL [https://transformer-circuits.pub/2024/](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html)
636 [scaling-monosemanticity/index.html](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html).

637

638 Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide
639 shut? exploring the visual shortcomings of multimodal llms. In *Proceedings of the IEEE/CVF*
640 *Conference on Computer Vision and Pattern Recognition*, pp. 9568–9578, 2024.

641

642 Constantin Venhoff, Ashkan Khakzar, Sonia Joseph, Philip Torr, and Neel Nanda. Too late to recall:
643 The two-hop problem in multimodal knowledge retrieval. In *Mechanistic Interpretability for Vision*
644 *(Non-proceedings Track)*, *CVPR 2025*, 2025a. URL [https://openreview.net/forum?](https://openreview.net/forum?id=VUhRdZp8ke)
645 [id=VUhRdZp8ke](https://openreview.net/forum?id=VUhRdZp8ke).

646

647 Constantin Venhoff, Ashkan Khakzar, Sonia Joseph, Philip Torr, and Neel Nanda. Too late to recall:
648 The two-hop problem in multimodal knowledge retrieval. In *CVPR 2025 Workshop on Mechanistic*
649 *Interpretability of Vision (MIV)*, 2025b. Non-proceedings Track Poster.

650

651 Constantin Venhoff, Ashkan Khakzar, Sonia Joseph, Philip Torr, and Neel Nanda. How visual repre-
652 sentations map to language feature space in multimodal llms. *arXiv preprint arXiv:2506.11976*,
653 2025c.

654

655 Guowei Xu, Peng Jin, Li Hao, Yibing Song, Lichao Sun, and Li Yuan. Llava-o1: Let vision language
656 models reason step-by-step. *arXiv preprint arXiv:2411.10440*, 2024.

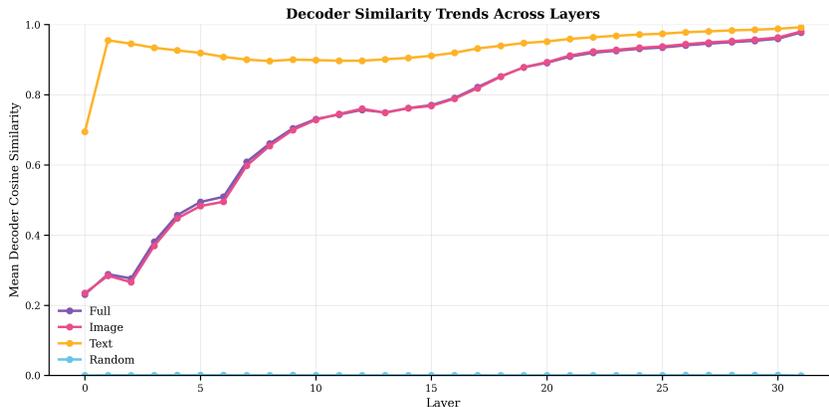
648 Fred Zhang and Neel Nanda. Towards best practices of activation patching in language models:
649 Metrics and methods. In *International Conference on Learning Representations (ICLR)*, 2024.
650 URL <https://doi.org/10.48550/arXiv.2309.16042>. arXiv:2309.16042.
651

652 Ruohong Zhang, Bowen Zhang, Yanghao Li, Haotian Zhang, Zhiqing Sun, Zhe Gan, Yinfei Yang,
653 Ruoming Pang, and Yiming Yang. Improve vision language model chain-of-thought reasoning.
654 *arXiv preprint arXiv:2410.16198*, 2024.
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

702 A APPENDIX

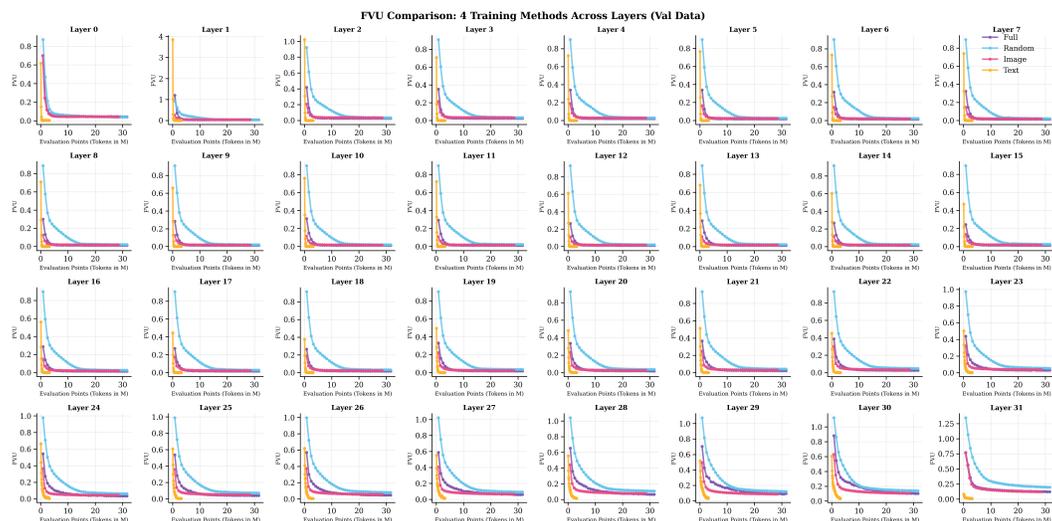
703 A.1 GEOMETRY DIVERGENCE: DECODER COSINE TRENDS

704 To quantify how SAE feature geometry shifts across training regimes, we track cosine similarity
 705 between decoder directions from SAEs trained on different input types. Fig 5 shows that text-only
 706 SAEs remain closely aligned across layers, while image-only and full-sequence SAEs diverge in early
 707 layers before realigning deeper in the model. Randomly initialized SAEs stay largely uncorrelated,
 708 confirming the stability of the observed trends.
 709
 710



711
 712
 713
 714
 715
 716
 717
 718
 719
 720
 721
 722
 723
 724
 725
 726 **Figure 5: Decoder cosine similarity vs. layer (LLM SAE vs. VLM SAE).** Text-only stays highly
 727 aligned across layers; image-only and full-sequence rotate in shallow layers and align later; random
 728 remains near zero. Higher cosine indicates closer alignment of SAE decoder directions.
 729

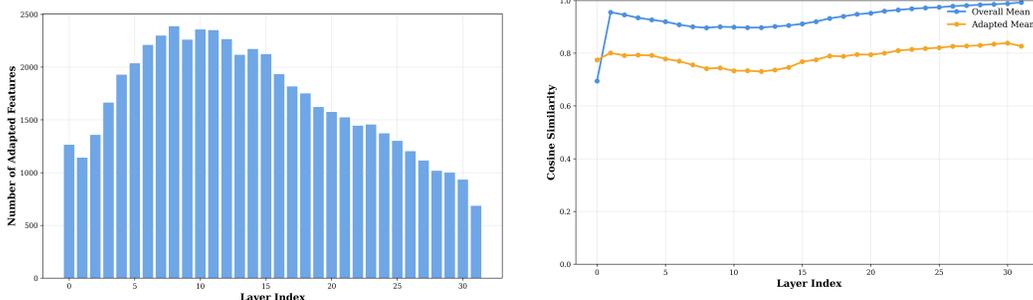
730 A.2 PER-LAYER FVU TRAJECTORIES



731
 732
 733
 734
 735
 736
 737
 738
 739
 740
 741
 742
 743
 744
 745
 746
 747
 748
 749
 750 **Figure 6: Per-layer FVU across regimes.** Each panel shows the convergence of SAEs trained with
 751 different masking regimes for a specific layer. Text-only SAEs begin with slightly higher error in
 752 the shallowest layers but adapt almost immediately to near-zero reconstruction. Image-only and
 753 full-sequence SAEs converge more slowly and plateau at higher error, while random initialization
 754 performs worst throughout. This confirms that projector-driven spans remain off-distribution in early
 755 layers and only align with the LLM basis in later layers.

A.3 PER-LAYER STATISTICS

Fig 7 shows that adapted features cluster in mid layers and taper in deeper blocks. Their decoder directions remain less aligned to the base dictionary than the overall pool, confirming stronger rotations under multimodal fine-tuning.



(a) **Adapted features per layer.** Most concentrate in mid layers, tapering in deeper blocks.

(b) **Decoder cosine by layer.** Adapted features remain less aligned to the base dictionary than the overall pool.

Figure 7: **Per-layer statistics of adapted features.** (a) Distribution of adapted feature counts across depth. (b) Mean decoder cosine similarity for adapted features vs. the overall pool.

A.4 THRESHOLD SWEEP FOR FEATURE SELECTION

To ensure that our choice of thresholds is robust, we sweep over the cosine percentile cutoff (p_{COS}) and visual energy threshold (ϵ). Fig. 8 reports three metrics: (i) total number of selected features, (ii) Jaccard overlap with the baseline adapted set, and (iii) per-layer count correlation. The results show a broad stable region around $\epsilon \approx 10^{-3}$ and $p_{\text{COS}} \approx 25\%$, which yields a compact yet consistent set of adapted features. We adopt this operating point (white circle) for all downstream analyses.

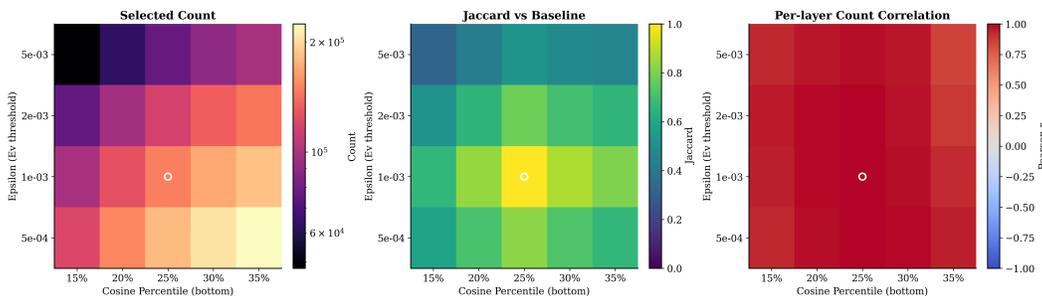


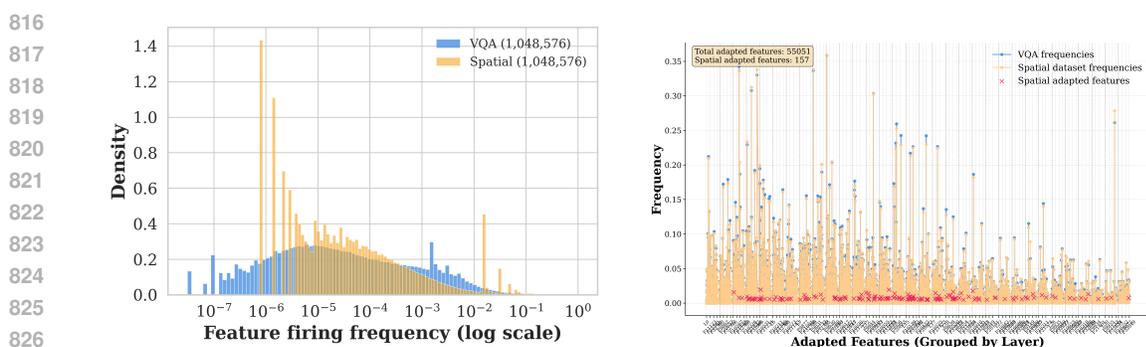
Figure 8: **Threshold sweep for feature selection.** Left: feature counts increase smoothly with more lenient thresholds. Middle: Jaccard overlap with the baseline peaks near the chosen point. Right: per-layer counts remain highly correlated across thresholds. The white circle marks the adopted operating point.

The visual-energy statistic E_v is computed under a text-only mask, since our SAEs are text-only. As a result, most features have $E_v = 0$, so requiring $\epsilon > 0$ acts as a strong filter. When cross-checking with downstream spatial tasks, we find that features with very low E_v rarely contribute meaningfully: they tend to cluster in shallow layers, show low spatial hit rates, and often appear polysemantic on inspection. In contrast, those that pass the ϵ cutoff carry a cleaner visual signal and align more consistently with spatially selective units in downstream evaluations, suggesting that the thresholded set captures genuinely vision-grounded features.

810
811

A.5 DISTRIBUTION-SHIFT VISUALIZATIONS

812 To complement the main-text description of our feature-selection procedure, we include here the
813 firing-frequency distributions and candidate-feature scatter plots used to identify spatial units under
814 different prompting conditions.
815



827 (a) Firing-frequency distributions for $\mathcal{D}_{\text{base}}$ and the
828 spatial split \mathcal{D}_{sp} .

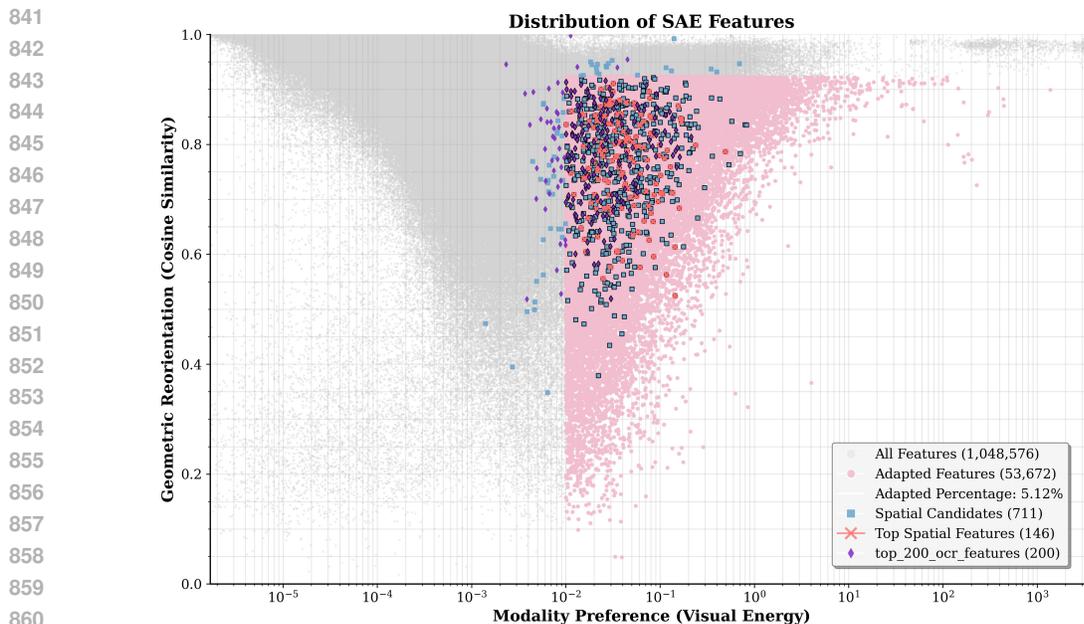
827 (b) Spatial candidate features under both splits, with
828 selected units highlighted.

829
830 **Figure 9: Spatial distribution shift.** Visualization of feature firing frequencies and candidate
831 selection under the spatial vs. base splits.
832

833

A.6 OCR FEATURE VISUALIZATIONS

834
835
836 We also apply our distribution-shift procedure to OCR-style prompts (e.g., “What does the sign
837 say?”). Fig. 10 shows that OCR-selective features cluster within the same adapted region as the
838 spatial subset, indicating that multimodal fine-tuning concentrates visually grounded capabilities into
839 a compact envelope of feature space.
840



861
862 **Figure 10: Distribution of OCR features.** Top OCR candidates (purple) cluster among adapted units
863 (pink), paralleling the spatial subset (blue).

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

B ADDITIONAL AUTO-INTERP EXAMPLES

In the main text (Sec. 5.1), we showed examples of adapted features using our automated interpretation pipeline. We include two further examples here. In both cases, the top-activating samples agree across VQA and VSR, and the interpretations are consistent and monosemantic.



Figure 11: **Additional Auto-Interp** examples. Top-activating VQA and VSR samples for two adapted features, showing consistent spatial relations.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

B.1 AUTO-INTERPRETATION AND SCORING PIPELINE

We evaluate interpretability using an automated feature-description pipeline with two variants: *RAW* (image+text) and *OVERLAY* (image+text+top-head heatmaps). For each feature f : 1. Select up to $k=5$ top-activating samples (deduped across VQA/VQA-spatial/VSR). 2. Call the API once to generate a single concise description. 3. Validate using held-out positive samples and random VQA negatives (two short rounds). 4. Compute F1 as a lightweight proxy for description confidence.

Outputs are stored per feature as JSON (`description`, `examples`, `classification results`). Adding overlays improves interpretability, with early results showing a typical gain of about +0.2 F1.

PROMPT A: Description (RAW / OVERLAY)

System: You are analyzing individual neurons using their top-activating samples (image+text; OVERLAY also includes attention heatmaps).

Task: Produce *one* short, lower-case sentence completing: “this neuron activates for ...”.

Guidelines: Base it on consistent patterns supported by image (+ overlays) and text; be specific; no hedging.

Return: { "description": "one concise sentence" }.

PROMPT B: Validation (F1)

System: You are validating a neuron description against short examples (image+text; OVERLAY adds heatmaps).

Task: For each sample, output 1 if it reasonably matches the description; else 0.

Return: { "classifications": [0/1, ...] }.

C ATTRIBUTION PATCHING ADDITIONAL EXPERIMENTS RESULTS

C.1 AGGREGATED ATTRIBUTION RESULTS

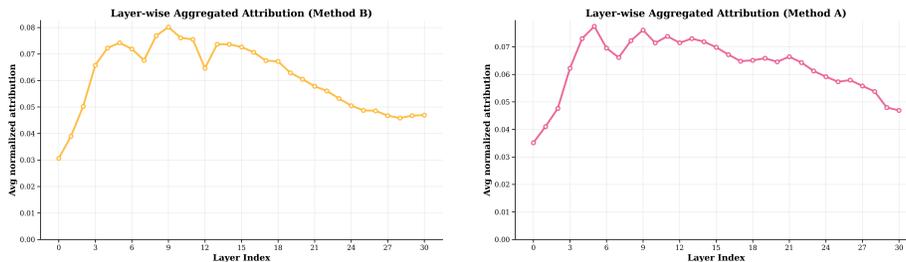


Figure 12: Layer-wise aggregated attribution curves for Method B (left) and Method A (right). Both peak in around middle layers, consistent with the emergence of spatial features.

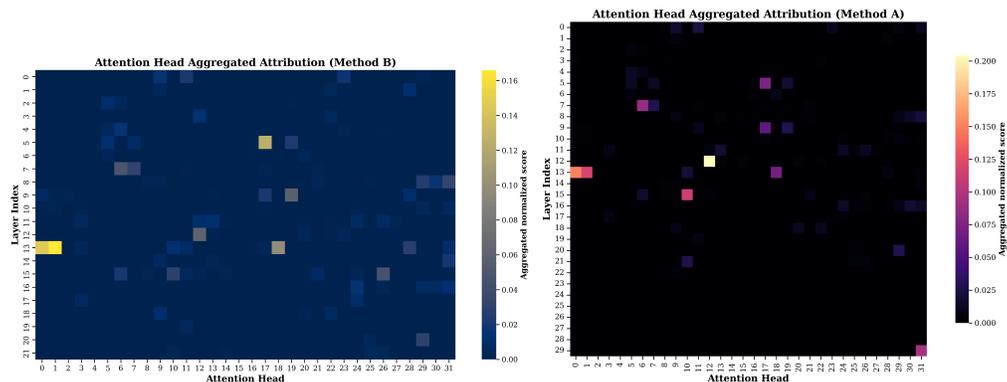
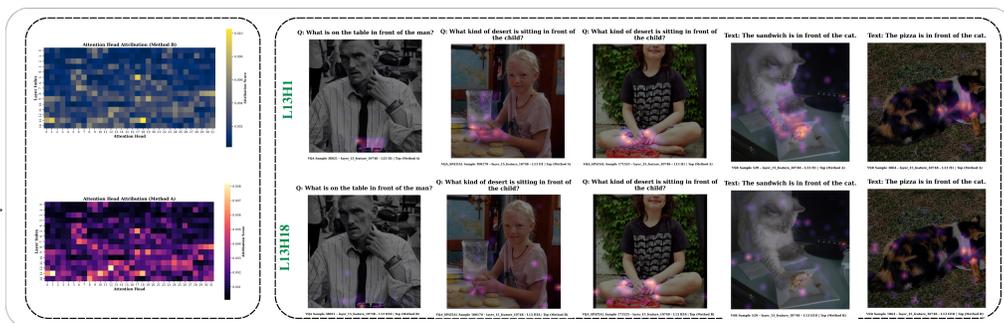


Figure 13: Attention head aggregated attribution maps for Method B (left) and Method A (right). Both highlight a similar set of specialized heads with high attribution scores.

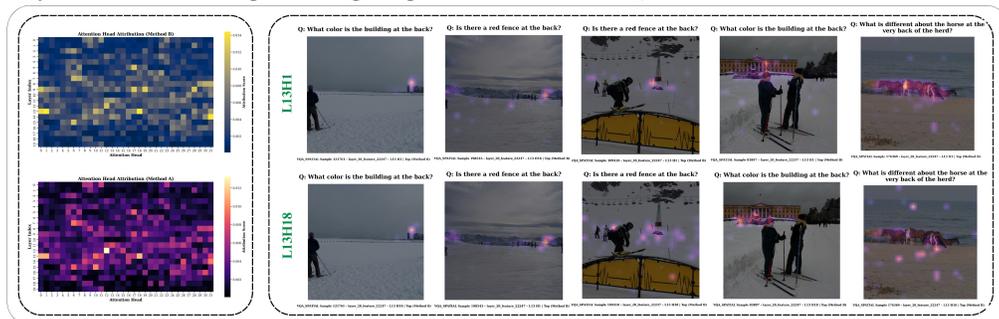
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

C.2 PER-FEATURE PANELS WITH TOP HEADS

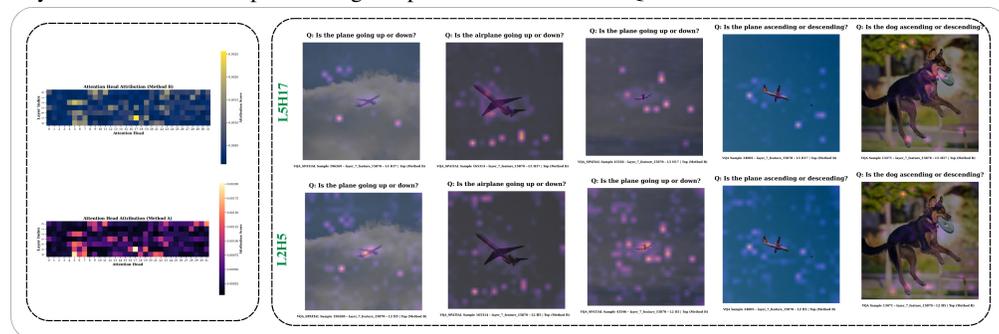
For individual spatial features, we show (i) per-layer/head attribution maps (Methods A and B) and (ii) attention overlays from the strongest heads on the feature’s top-activating samples across both VSR and VQA datasets.



(a) **Layer 15, Feature 10748.** VSR Relation: “in front of.” Top heads (Method A): L13H1, L12H12, L13H18. Top heads (Method B): L13H18, L5H17, L13H1. *Overlap:* L13H1, L13H18. Attention overlays are shown on the top-activating samples across VSR and VQA.



(b) **Layer 20, Feature 22247.** VSR Relation: “at the back of.” Top heads (Method A): L12H12, L13H18, L13H1. Top heads (Method B): L13H1, L13H18, L14H31. *Overlap:* L13H1, L13H18. Attention overlays are shown on the top-activating samples across VSR and VQA.



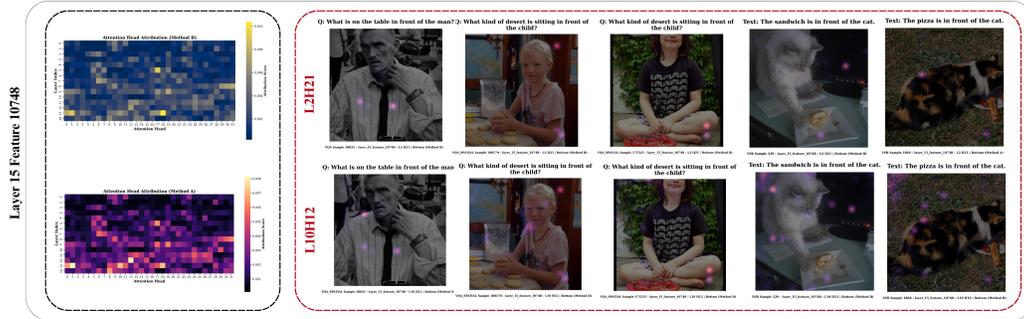
(c) **Layer 7, Feature 15870.** VSR Relation: “above.” Top heads (Method A): L5H17, L6H5, L0H31. Top heads (Method B): L5H17, L2H5, L2H6. *Overlap:* L5H17. Attention overlays are shown on the top-activating samples across VSR and VQA.

Figure 14: **Attribution patching on individual spatial features.** Each subfigure displays aggregated head/layer attribution maps (left) and attention overlays (right) using the strongest heads on the feature’s top-activating samples across both VSR and VQA.

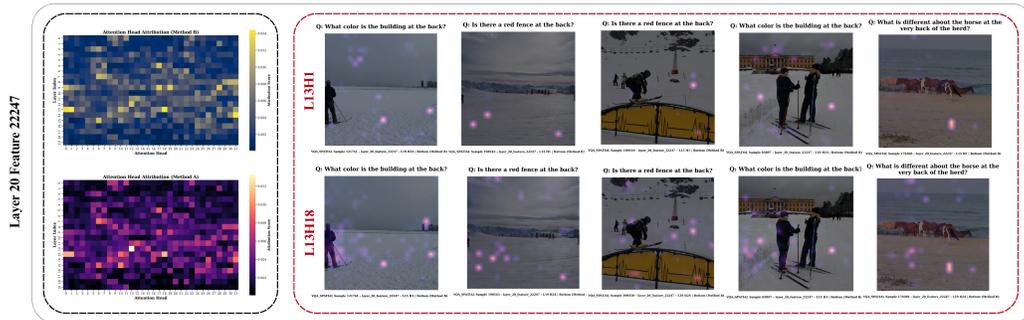
Across these examples, the two attribution methods consistently surface overlapping heads, indicating that a small group concentrates much of the spatial signal. Method B generally produces sharper rankings and cleaner overlays, suggesting it is more reliable for identifying the causal drivers of spatial features.

C.3 BOTTOM-RANKED HEADS AS A CONTROL

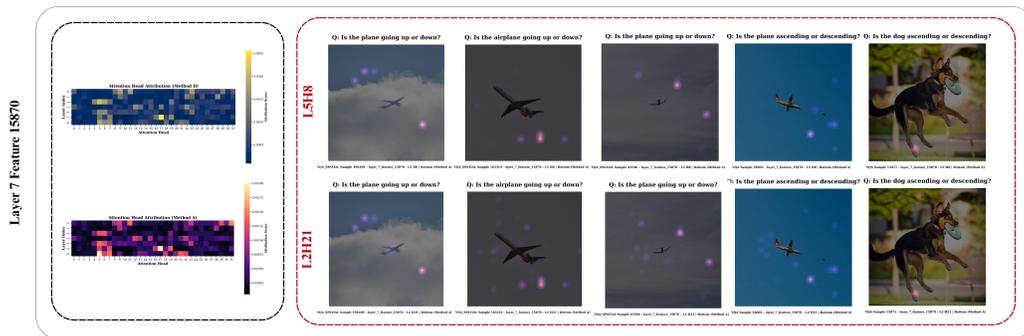
As a control, we visualize overlays from the *bottom-ranked* heads (per method, per feature). Across VSR and VQA top-activating samples, these heads generally fail to localize semantically relevant regions.



(a) Layer 15, Feature 10748. VSR Relation: “in front of.”



(b) Layer 20, Feature 22247. VSR Relation: “at the back of.”



(c) Layer 7, Feature 15870. VSR Relation: “above.”

Figure 15: **Bottom-ranked heads yield weak localization.** For each feature, we show overlays from the lowest-scoring heads under Methods A and B on the feature’s top-activating samples across VSR and VQA. In contrast to Appx. Fig. 14, these heads produce diffuse or irrelevant attention.

D FULL ABLATION RESULTS

Table 2 reports a more detailed version of ablation results for the top SAE features. For each feature, we show average accuracy and probability drops on VSR across seeds, together with the number of evaluation samples. We also report accuracy drops on VQA, random-feature control drops (Δ Ctrl), odds ratios (VSR OR), and relation-specific subsets of VSR derived from top-activating samples.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

Large negative Δ VSR Acc with small Δ VQA Acc indicates spatial specificity, near-zero Δ Ctrl supports robustness, and high odds ratios reflect selective recruitment under spatial prompts.

Layer	Feature	Δ VSR Acc	Δ VSR Prob	Δ VQA Acc	Δ Ctrl	VSR OR	#Samples	VSR Relations
11	27061	-13.30	-0.09	-0.40	0.00	8.03	94	across from
12	23874	-10.24	-0.10	-0.40	-0.95	9.10	210	left of
18	29948	-7.98	-0.09	-0.30	0.00	8.36	188	beside
23	4060	-5.85	-0.00	-0.70	-1.06	7.47	94	at the back of
14	17873	-10.00	-0.07	-0.30	-2.71	7.17	480	at the right side of
9	15404	-11.19	-0.07	-0.80	1.08	5.60	277	below
7	6986	-10.87	-0.03	-0.50	0.34	4.74	589	under
7	15870	-15.54	-0.09	-0.10	-0.88	4.32	341	above
10	5121	-7.92	-0.06	-0.10	0.12	4.22	846	above, on top of
11	24089	-7.68	-0.05	-0.60	-0.12	4.18	846	above, on top of
12	13305	-6.38	-0.05	-0.70	0.24	4.17	846	above, on top of

Table 2: Full ablation results for top SAE features, averaged over seeds. The number of VSR samples evaluated is shown alongside accuracy/probability drops and odds ratios.

E

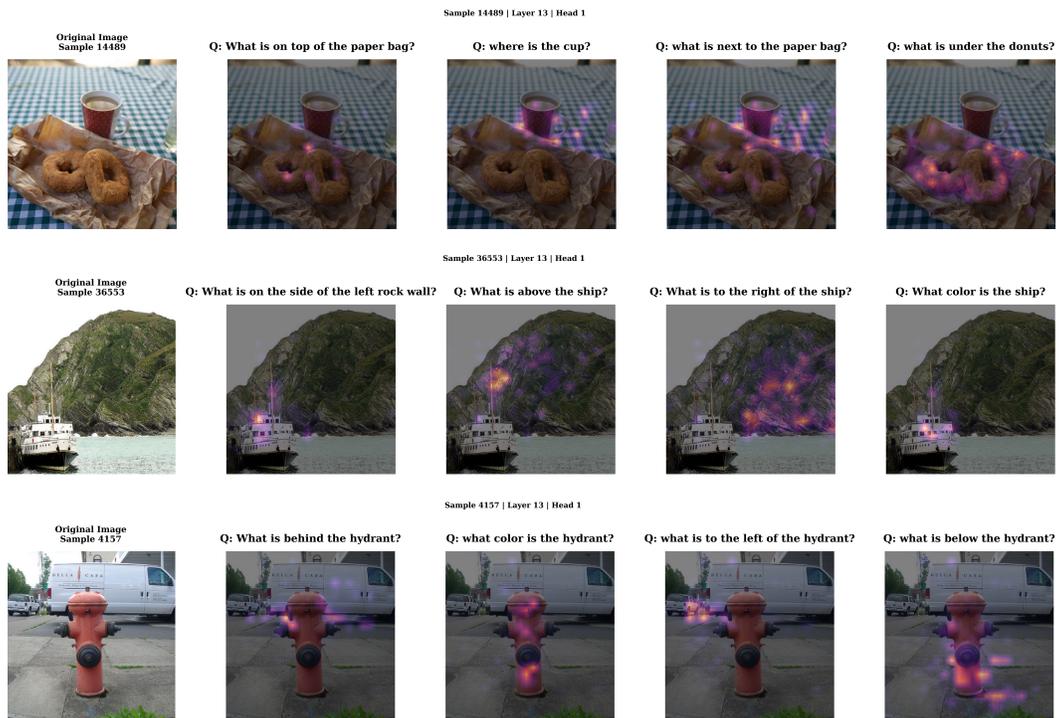


Figure 16: **Attention head visualizations across queries.** Each row shows one image with attention overlays from a single high-attribution head across multiple spatial and non-spatial custom queries. The same heads consistently focus on semantically relevant regions.