

LINEAR MECHANISMS FOR SPATIOTEMPORAL REASONING IN VISION LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Spatio-temporal reasoning is a remarkable capability of Vision Language Models (VLMs), but the underlying mechanisms of such abilities remain largely opaque. We postulate that visual/geometrical and textual representations of spatial structure must be combined at some point in VLM computations. We search for such confluence, and ask whether the identified representation can causally explain aspects of input-output model behavior through a linear model. We show empirically that VLMs encode object locations by linearly binding *spatial IDs* to textual activations, then perform reasoning via language tokens. Through rigorous causal interventions we demonstrate that these IDs, which are ubiquitous across the model, can systematically mediate model beliefs at intermediate VLM layers. Additionally, we find that spatial IDs serve as a diagnostic tool for identifying limitations and bottlenecks in existing VLMs. We extend our analysis to video VLMs and identify an analogous linear temporal ID mechanism. By characterizing our proposed spatiotemporal ID mechanism, we elucidate a previously underexplored internal reasoning process in VLMs, toward improved interpretability and the principled design of more aligned and capable models.

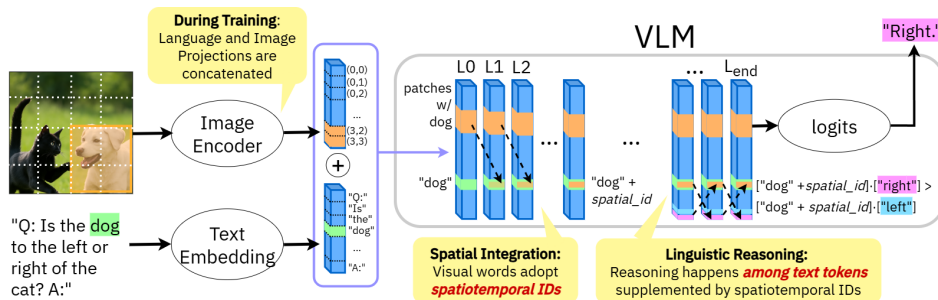


Figure 1: **Hypothesis for spatiotemporal visual reasoning.** The VLM linearly binds spatiotemporal localization to object word activations in early layers. Subsequent linguistic reasoning about the object is informed by its location in space and time per the spatiotemporal ID.

1 INTRODUCTION

Reasoning about visual input with textual queries is a key challenge behind vision-language models (VLMs). Consider a typical visual question answering (VQA) prompt: “*Is the dog to the left or right of the cat?*”. To succeed at this, one must resolve linguistic references, locate entities in the visual field, assess spatial relationships, and make a categorical decision. Though complex capabilities in spatial or temporal reasoning are still far from being fully understood or reliably engineered (Stogiannidis et al., 2025; Chen et al., 2025; Tong et al., 2024), SoTA VLMs have seen steady progress in simple visual reasoning without explicit guidance. So how do they do it?

Attention-based analyses in VLMs have shown various structural properties emerge in VLM internals during VQA (Jiang et al., 2025b; Neo et al., 2024; Zhang et al., 2024a). Relatedly, mechanistic interpretability in LLMs has uncovered linear circuits for relational linguistic reasoning (Park et al., 2024; Feng & Steinhardt, 2024; Hernandez et al., 2024). Might such linear processes also be driving visual reasoning in VLMs, and if so, how exactly? This leads us to ask: **Q1.** *Can we linearly model emergent structured reasoning processes that drive spatial reasoning in VLM internals?*

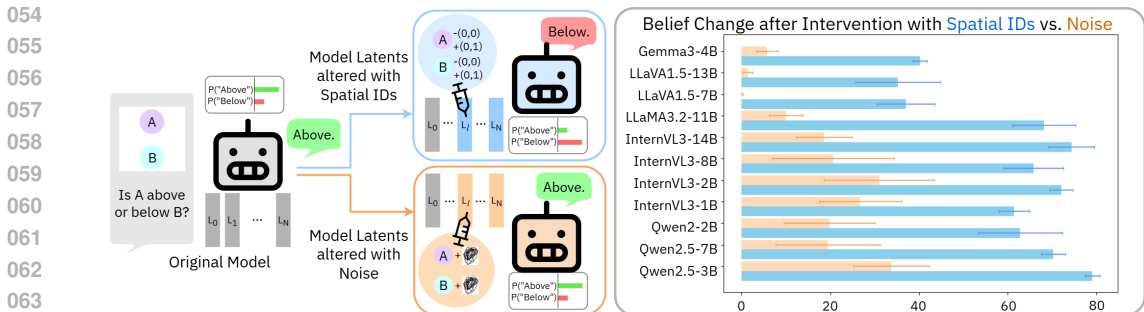


Figure 2: **Results from Targeted Intervention** (§3). Median binary belief swap due to spatial ID steering is 64.4%, and 29.5% for noise. Spatial IDs have 43.6% above-chance influence on average. We conclude that spatial IDs mediate models’ beliefs about objects’ locations in space.

The typical VLM architecture utilizes a vision encoder which projects the input image to embeddings that are prepended to text token embeddings. This is then processed by a downstream vision-aligned LLM. Popular model families using this paradigm are LLaVA(Liu et al., 2023), LLaMA(Dubey et al., 2024), Qwen (Bai et al., 2025), InternVL (Chen et al., 2024b), and Gemma (Team et al., 2024). A growing body of work aims to improve spatial reasoning capacities in VLMs (Chen et al., 2024a; Fan et al., 2025) and temporal reasoning in video models (Xiao et al., 2024; Li et al., 2024b). Identification of the internal mechanism by which SoTA VLMs do spatial VQA can empower engineers to identify current architectural components leading to VQA failure modes in 3D reasoning or simple VQA. To this end, we ask: **Q2. Given our linear model of spatial reasoning in model activations, how do we use it to understand and improve SoTA VLMs?**

Similar training paradigms to image-based VLMs yield video models such as LLaVA-Video(Zhang et al., 2024b), VideoLLaMA3(Zhang et al., 2025), and Qwen2.5 (Bai et al., 2025), among others. Given our theory for the mechanisms underlying spatial reasoning in VLMs, we ask: **Q3. Do video models utilize analogous linear mechanisms for temporal reasoning?**

To address these questions, we conduct a mechanistic analysis of autoregressive VLMs and construct a linear model for spatiotemporal reasoning. We show that VLMs decompose a visual reasoning task by first binding spatial information about visual objects to object word activations, in the form of linear components we term *spatial IDs*, answering Q1 (Fig. 1). We then extract these IDs and demonstrate their mediative capacity on model output through targeted belief steering in text activations (Fig. 2). We further find that spatial IDs provide insight on VLMs’ struggle with depth reasoning, and incorrect spatial IDs as a result of weak vision encoder or poor modality integration leads to failures in LLaVA and LLaMA. This answers Q2. Finally, we show that temporal IDs similarly mediate video models, answering Q3. In summary, our novel contributions are:

- **Spatial ID Model Formulation:** We propose a linear model of spatial reasoning in VLMs, called *spatial IDs*. These are text-anchored latent structures that bind visual elements to object tokens thus enabling linguistic reasoning about space (§2.1). We empirically extract them from SoTA VLMs for characterization (§2.2).
- **Analytical and Empirical Proof of Causality:** We show model belief can be manipulated by perturbing only the spatial IDs, demonstrating their causal role in reasoning (§3), and provide theoretical intuition for the emergence of spatial IDs in VLMs (§2.3).
- **SoTA VLM Analysis with Spatial IDs:** Through targeted intervention, we identify limitations in depth expression (§4.1) and systematic failure modes in LLaMA/LLaVA (§4.2), and find models can be effectively finetuned with spatial ID guidance (§4.3).
- **Extension to Temporal IDs in Video Models:** We perform our extraction and characterization analysis on SoTA video models and show that linear temporal IDs, like spatial IDs, can drive temporal reasoning in VLMs (§5).

2 EMERGENT STRUCTURE IN SPATIAL VISUAL REASONING

In this section, we characterize the spatial reasoning circuits in SoTA VLMs and isolate any linearly separable components used to communicate spatial information. Towards this end, we track information flow in VLMs and identify important junctions for spatial information transfer across token sequences. Then we empirically extract linear spatial IDs, and analytically derive how they arise.

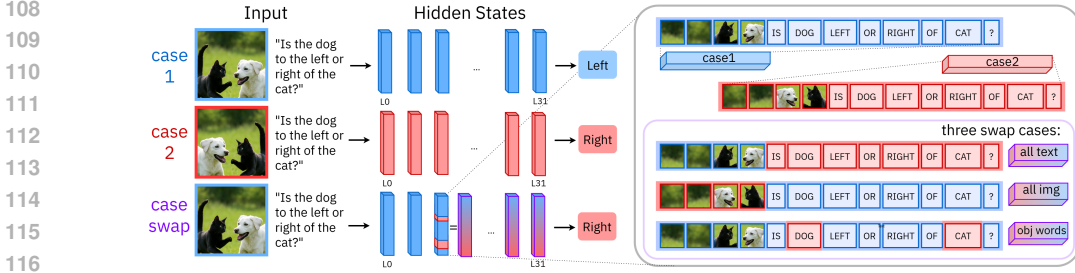


Figure 3: **Mirror swapping experiment** (§2.1). Activations from case 1 and 2 are partially swapped at a select layer, in one of three arrangements. Computations continue normally after this point.

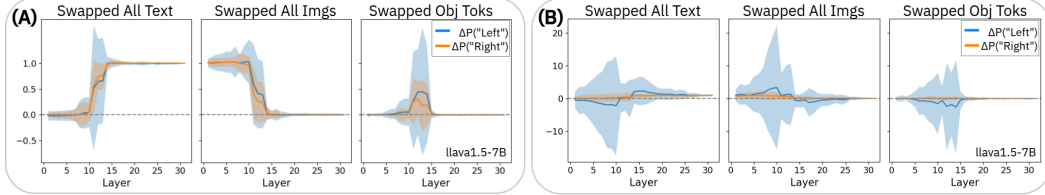


Figure 4: **Ratio change in log probability for logits “left” and “right” from mirror swap (A) and attribute swap (B) interventions.** (A) shows distinct binary belief swaps, where text tokens have an influence after middle layers. Image patches stop having an influence after that point, and object word tokens *only* have an influence in these middle layers. The control, (B), is noisy.

2.1 TRACKING INFORMATION FLOW DURING REASONING

To uncover whether VLMs engage in structured visual reasoning, i.e., isolating and propagating spatial information across layers, we intervene on internal activations during inference.

Mirror Swapping Experiment. Our goal is to compare the model’s output when presented with two distinct images and the same text query. If the model uses localized intermediate representations to reason about spatial relationships, then swapping activations between spatially distinct inputs at key layers and sequence indices should disrupt its final belief about spatial orientation, while swaps between spatially equal but attribute-wise different inputs shouldn’t have a strong effect.

Concretely, we run inference on plain and mirrored image-text pairs, extract their representations x at an intermediate layer L , then replace a subset Q of activations in the original x_L with activations from the mirrored counterpart y_L . The modified representation \tilde{x}_L is passed through the remaining layers. We conduct interventions with three variants of Q : (1) all text tokens (2) all image patches (3) object-word tokens only. If information critical to spatial reasoning is concentrated in any of these, the model’s belief will change when that region is overwritten. As a control, we concurrently perform “attribute swapping”, which follows the same steps but instead of mirroring the input image for the intervention case, changes its colors. The intervention procedure is visualized in Fig. 3 and formally defined in Alg. 1. Further implementation specifics are deferred to Appendix §A.1.

Algorithm 1 Swapping Intermediate Activations in Mirrored Images

$$\begin{aligned}
 x_L, y_L &\leftarrow f_L \circ \dots \circ f_1(x), \quad f_L \circ \dots \circ f_1(y) && \triangleright x, y: [\text{seq_dim}, \text{embed_dim}] \\
 \tilde{x}_L &\leftarrow x_L[\tilde{Q}] + y_L[Q] && \triangleright \tilde{x}_L: [\text{seq_dim}, \text{embed_dim}], Q: [\text{num_of_inds}] \\
 \tilde{x}_{\text{out},L}, y_{\text{out}} &\leftarrow f_{L_{\text{max}}} \circ \dots \circ f_{L+1}(\tilde{x}_L), \quad f_{L_{\text{max}}} \circ \dots \circ f_{L+1}(y_L) && \triangleright P_{\tilde{x}_{\text{out},L}}(\text{“GT”}): [1]
 \end{aligned}$$

Here, Q denotes the selected indices in the input sequence to swap, and \tilde{Q} is all other indices. We use the COCO-SPATIAL benchmark (Kamath et al., 2023) for the mirrored images, which is a curated subset of COCO (Lin et al., 2014) annotated with spatial language. To quantify belief shift caused by the intervention, we compute the fraction of the mirror-induced change that can be attributed to the swapped activations at layer L . For the ground truth logit “GT”, this quantity is derived as:

$$\text{belief shift}_L = \frac{P_{x_{\text{out},L}}(\text{“GT”}) - P_{\tilde{x}_{\text{out},L}}(\text{“GT”})}{P_{x_{\text{out},L}}(\text{“GT”}) - P_{y_{\text{out}}}(\text{“GT”})} \quad (1)$$

Results from Mirror Swapping are shown in Fig. 4A. Through mirror swapping, we observe a *layer-specific effect* for intervention effect across modalities. Intervening on visual patch tokens

has a strong effect in early layers but fades with depth. Conversely, interventions on text tokens increasingly affect final outputs in later layers. This is corroborated by observations that middle layers have a modality switching effect in VLMs (Jiang et al., 2025b). Notably, swapping only the object-word tokens alters spatial belief specifically within a narrow band of intermediate layers.

Attribute swapping results (Fig. 4B) indicate that mirror swapping is a strong experimental setup for assessing spatial information flow in isolation from spurious visual factors. For the belief shift metric, a value of 0.0 on the y axis indicates model belief in the intervened case is equivalent to case 1 (original query), while 1.0 indicates the belief is equivalent to case 2 (mirrored/changed query). Mirror swapping results in distinct and strong binary belief swaps whereas attribute swapping yields mostly noise, to the point belief shift magnitudes are $-20\sim 20x$ that of the original belief difference.

These results suggest that VLMs extract and encode spatial facts from the image into object word tokens’ activations, then operate over them in text-space. We term the latent structures holding visual spatial information *spatial ids*. Inspired by latest mechanistic interpretability findings (discussed in §6), we hypothesize that the manner of spatial information storage is approximately linear.

2.2 EMPIRICAL DERIVATION OF SPATIAL IDS

If spatial IDs are indeed linearly bound to object word activations, we should be able to extract them by averaging out object-related semantics from text activations. Below we outline the process of their extraction. In §3, we will test if these IDs causally mediate model beliefs, to validate whether the spatial reasoning mechanism in VLMs is indeed linear.

Extraction Preliminaries. We first set up some formalisms to derive spatial IDs. Let $\mathcal{O} = \{o_1, o_2, \dots, o_N\}$ denote a set of object categories. For each object $o \in \mathcal{O}$, we have a set of images $\{I_{(i,j)}\}$ where the object is positioned at spatial coordinates (i, j) in a $m \times m$ grid. Then let $T^{(o)}$ be a natural language query containing the token corresponding to object o , such as “Is there an o ?”. We define $\phi_L(o; I_{(i,j)}^{(o)}, T^{(o)}) \in \mathbb{R}^d$ as the embedding of the text token corresponding to object o , extracted from layer L of the VLM when input= $(I_{(i,j)}^{(o)}, T^{(o)})$. The mean embedding for object o at layer L is then:

$$\bar{\phi}_L^{(o)} = \frac{1}{m^2} \sum_{i=0}^{m-1} \sum_{j=0}^{m-1} \phi_L(o; I_{(i,j)}^{(o)}, T^{(o)}) \quad (2)$$

Yielding the object-specific spatial ID at location (i, j) for object o :

$$\Delta_L^{(o)}(i, j) = \phi_L(o; I_{(i,j)}^{(o)}, T^{(o)}) - \bar{\phi}_L^{(o)} \quad (3)$$

From this we can derive the *universal spatial ID* at location (i, j) , averaged over N objects.

$$\Delta_L(i, j) = \frac{1}{N} \sum_{n=1}^N \Delta_L^{(o_n)}(i, j) \quad (4)$$

To extract canonical horizontal and vertical directions from the universal spatial IDs $\Delta_L(i, j) \in \mathbb{R}^d$, we compute average difference vectors across grid-aligned coordinate pairs. The vertical and horizontal direction vectors $v_L, h_L \in \mathbb{R}^d$, corresponding to increasing i and j , are computed based on the spatial IDs. Eq. 5 shows the derivation for v_L , and h_L is derived in an analogous manner.

$$v_L = \frac{1}{m \cdot \binom{m}{2}} \sum_{i=0}^{m-1} \sum_{j_1 > j_2} [\Delta_L(i, j_1) - \Delta_L(i, j_2)] \quad (5)$$

Empirical Extraction. For our study, we extract spatial IDs from 11 SoTA VLMs, with synthetic images created from open-source OBJAVERSE (Deitke et al., 2023) objects. The object renders are paired onto a grid of $m = 4$ on top of random natural backgrounds. We provide further extraction details in Appendix §A.2, along with ablations showing extracted spatial IDs are invariant to chosen images §D and counterfactual studies confirming that spatial IDs reside in object words, and spatial

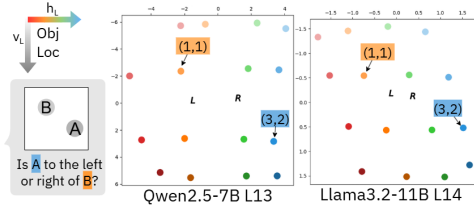


Figure 5: **Spatial IDs in a grid.** Color and saturation of markers represent the location of the object when spatial ID was extracted. x and y axes are coefficients of ID projections onto h_L and v_L . L, R represent “left”, “right” textual activations.

axes are orthogonal §C. Fig. 5 shows two example spatial ID grids projected onto their respective spatial vectors. IDs from more models are shown in §B. We see that these extracted IDs arrange in an approximate $m \times m$ grid at modality binding layers. Also projected are activations for spatial words, where we find that “left” is closer to leftmost spatial IDs, and “right” vice versa.

2.3 THEORETICAL SKETCH OF SPATIAL IDS

We now offer a quick, highly minimal analytical intuition for how the emergence of spatial IDs can be ubiquitous across many different models. Let $p = (i, j)$ be some coordinate on a $m \times m$ grid. Then for some query to a VLM, let the input sequence contain projected visual tokens $\{x_p\}$ for all p , and the query text tokens include an object token o . The residual update to o by one head is:

$$r_o \leftarrow r_o + W_{\text{out}} \sum_{p \in P} \alpha_{o \leftarrow p} v_p, \quad \alpha_{o \leftarrow p} \propto \exp\left(\frac{q_o^\top k_p}{\sqrt{d}}\right), \quad v_p = W_V x_p. \quad (6)$$

With cross-modal alignment, attention peaks at the true object patch p^* , giving $\delta r_o \approx W_{\text{out}} W_V x_{p^*}$. Decompose each patch as $x_p = s_p + P \psi(p) + \varepsilon_p$, where s_p encodes content, $\psi(p) \in \mathbb{R}^{d_\psi}$ is a shared positional basis (e.g. RoPE or learned 2D embeddings), P maps positional features into model space, and ε_p is small. We can now substitute $\phi_L(o; I_{p^*}, T^{(o)}) = r_o + \delta r_{o, p^*}$ into Eq. 3. A detailed derivation is in §2.2, but in summary we get:

$$\Delta_L(p^*) = \Delta_L(i, j) \approx \underbrace{W_{\text{out}} W_V P}_M \left(\psi(i, j) - \frac{1}{m^2} \sum_p \psi(p) \right). \quad (7)$$

M (fixed per model)

Thus, spatial IDs are approximately a linear transformation of a universal positional basis written into the object token by attention. Spatial logits are thus approximately linear readouts:

$$\ell(\text{LEFT}) - \ell(\text{RIGHT}) \approx (w_{\text{LEFT}} - w_{\text{RIGHT}})^\top \Delta_L(i, j), \quad (8)$$

so if $(w_{\text{LEFT}} - w_{\text{RIGHT}})^\top M$ aligns with the x -coordinate in ψ , the model prefers “left.” Empirically, a low-rank linear fit from positional encodings ψ to spatial IDs Δ_L explains most variance (e.g. rank-3 gives $R^2 \gtrsim 0.85$, see §E.2, Table 1). A more detailed derivation for $\Delta_L(i, j)$ for the multihead case is shown in Appendix §E.1. This is a particularly simplified setting, and real reasoning circuits in VLMs will involve a lot more noise and nonlinearities. The main takeaway is that VLM designs like Fig. 1 encourage models to endow text tokens with visual information, followed by linguistic reasoning. This information transfer, in its most simplified linear form, is in the form of spatial IDs.

In practice, the finegrained circuit employed by VLMs may be much more varied, distributed, and nonlinear. The spatial ID framework could capture just one component of a more complex system. But per Ockham’s Razor, spatial IDs are powerful due to their simplicity. In following sections, we demonstrate the mediative influence of this simple spatial ID model on final VLM outputs, and further show how spatial IDs can be leveraged to improve existing models and build stronger ones.

3 SPATIAL IDS MEDIATE MODEL BELIEFS

If spatial IDs capture the causal mechanisms behind spatial reasoning, we should be able to linearly subtract or add arbitrary IDs to object word activations and change the model’s belief about object location. In this section, we design and perform experiments on real naturalistic images to test that empirically derived spatiotemporal IDs have causal effects on model outputs on spatial VQA.

Steering with Arbitrary IDs Experiment. For some layer L , we denote the model residuals corresponding to the entire input sequence after that layer as x_L , and perturb its token activation at some index q to observe any effects on the output belief. Alg. 2 illustrates the process.

Algorithm 2 Intervention at Layer L via Residual Modification

$$\begin{array}{ll} x_L \leftarrow f_L \circ \dots \circ f_1(x) & \triangleright x: [\text{seq_dim}, \text{embed_dim}] \\ \tilde{x}_L \leftarrow x_L[:q] + [x_L[q] + \Delta_L(i, j) - \tilde{\Delta}_L(i, j)] + x_L[q+1:] & \triangleright \Delta_L(i, j): [\text{embed_dim}] \\ \tilde{x}_{\text{out}} \leftarrow f_{L_{\text{max}}} \circ \dots \circ f_{L+1}(\tilde{x}_L) & \triangleright P_{\tilde{x}_{\text{out}}}(\text{“GT”}): [1] \end{array}$$

Here we scale the norm of $\Delta_L(i, j)$ to be $\alpha |x_L[q]|$, and $\tilde{\Delta}_L(i, j) = \Delta_L(m - i - 1, j)$. This approximately preserves the norm of x_L . $\alpha = 5$ is some scaling constant set after grid searching for

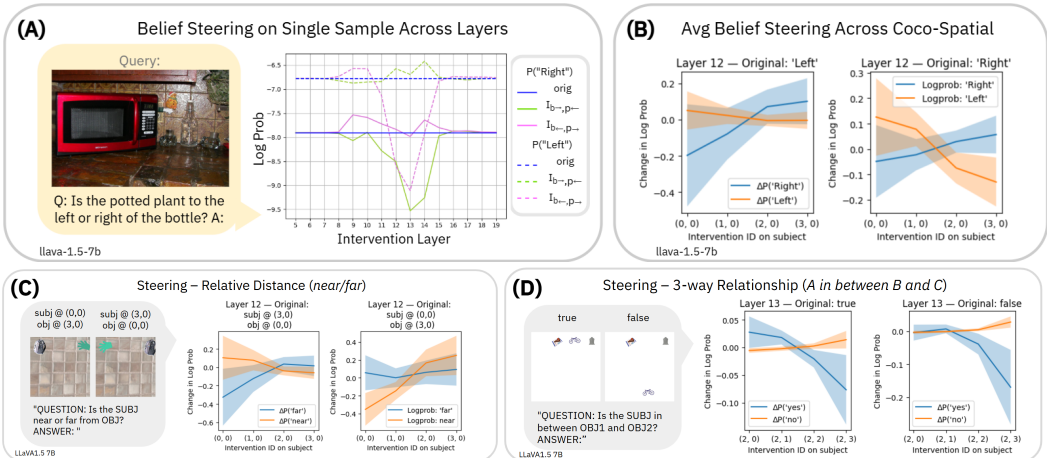


Figure 6: **Effect of spatial steering on real images** on one sample across different intervention layers (A) and across a dataset for one layer (B). In (A), dotted and solid lines indicate answer probabilities for “left” or “right”. Different colors indicate no intervention (blue), steering the bottle to the left and plant to the right (pink), and the reverse steering (green). Blue lines are flat and show that the unintervened model incorrectly assigns a higher log probability to “left”. Pink lines show intervention on intermediate layers results in overwriting initial incorrect beliefs. (B) shows the shift in log probability for “left” vs. “right” as a result of spatial steering on the subject word token. (C) and (D) show shifts in log probability for “near” vs. “far”, and “yes” to an object being sandwiched between two others, vs. “no”.

stable intervention. We take 100 COCO images where one object is to the left or right of another, per labels from COCO-SPATIAL, and ask queries of the form “Is x to the left/right of y?”. We measure the log probability of “left” and “right” tokens in the final output logits to assess steering effects.

Results from Arbitrary Steering. Fig. 6 shows the effects of model belief steering on real images and videos. Fig. 6A shows that steering impact is greatest at modality alignment layers as expected per the mirror swapping analysis, and Fig 6B shows that intervening with the rightmost spatial ID largely enhances model belief that the object is to the right, and vice versa for the leftmost ID for leftward belief. The y axes show changes in log probability for those binary directions for the whole dataset, and x axes show the different ID locations. Regardless of whether the answer to the original query was “left” or “right”, subplot trends are the same.

We repeat the analysis for queries about relative distance and three-way relationships where one object is sandwiched *in between* two others. Again, we find that when the object is to the left, altering the spatial ID of the subject towards the right increases the likelihood of “far” and decreases that of “near”, and vice versa if the object is to the right. Similarly, we find that bringing a subject closer and closer to be surrounded by two objects increases the model’s belief that the subject is *in between* the objects.

Adversarial Steering Experiment. If spatial IDs are indeed ubiquitous across models, interventions on internal activations should change the resultant model beliefs across many SoTA models. To confirm this, we evaluate the log probability of the correct answer (“GT”) and its opposite (“-GT”) for all samples in COCO-SPATIAL on 11 SoTA models. Then, we repeat this measurement after intervention with spatial IDs most likely to reverse their original beliefs. More detailed experimental procedure is provided in §A.5. In addition to targeted adversarial steering, we perform steering with noise vectors of the same norm as the spatial IDs, to evaluate chance belief swaps.

Adversarial Steering Results. We report % binary belief swaps on COCO-SPATIAL from the spatial ID vs. noise steering case in Fig. 2. Steering with spatial IDs yields a median 64.6% swap in beliefs, versus 29.5% with noise. This indicates activation intervention has nonzero chance influence on model output, but there is a clear above-chance average of 43.6% increase with spatial IDs. Here, a model’s belief on one sample is considered “swapped” if the relative likelihood of the ground truth and its opposite answer has changed. For example, if $P(\text{“left”}) > P(\text{“right”})$ before intervention, but after intervention we see $P(\text{“left”}) < P(\text{“right”})$, the intervention has swapped the model belief. Thus we conclude that spatial ID mechanisms mediate model belief in the models considered.

4 SPATIAL IDS FOR UNDERSTANDING AND IMPROVING IMAGE VLMS

With the existence and causal nature of spatial IDs established, we explore two ways to leverage them towards stronger VLMS. First, we aim to understand why 3D reasoning fails in SoTA VLMS. Second, we use spatial IDs to diagnose architectural bottlenecks of SoTA VLMS in VQA.

4.1 DEPTH REPRESENTATION IN IMAGE VLMS

Spatial IDs suggest that VLMS represent visual space within a 2D grid. What might this mean for depth? We hypothesize that the language model must reason about depth related queries using the 2D localization in context. To verify whether this is the case, we look at the resulting belief changes in the depth axis when the LLaVA1.5 7B model is steered with spatial IDs varying in height. Fig. 7 shows the results. The same spatial IDs increasing the likelihood for “above” and decreasing “below”, also drive up “front” and drive down “behind” in LLaVA.

Further, projection of these word embeddings onto spatial vectors reveals that “above”/“behind” and “below”/“front” map to overlapping locations, indicating their functional relationships with spatial IDs are similar. These results may be due to biases in training, or innate shortcomings in the VLM architecture. They certainly highlight the need for better depth-handling mechanisms, whether that be through improved training data or tooling.

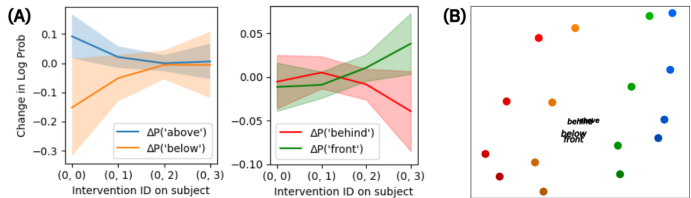


Figure 7: **Depth and height are strongly correlated in LLaVA.** (A) Steering results for IDs varying in y-dim and their impact on beliefs about height or depth. (B) Projection of spatial words onto a spatial ID grid. Embeddings for “above”/“front” and “behind”/“below” are nearly identical.

4.2 DIAGNOSING VLMS

When a VLM fails at a spatial task, how do we pinpoint the reason it failed? Referring back to Fig. 1, VLM failure points can roughly be divided into modality encoding, crossmodal information integration, or linguistic reasoning stages. Knowing what part of a VLM’s architecture must be improved to reduce failures is paramount to efficient model engineering.

Per-sample analysis of spatial IDs provides a unique ability to identify a model’s bottleneck. Consider an evaluation set $\mathbf{K} = \{k_1, k_2 \dots k_K\}$, where each $k = (image, query)$. An imperfect VLM will fail at some samples. In this section, we perform two experiments to identify the architectural component which causes for the distribution of \mathbf{K}_{wrong} to be statistically distinct from $\mathbf{K}_{correct}$.

An example diagnosis process may look like this. If a model exhibits *incorrect* spatial ID binding, and that incorrect output produced is faithful with the spatial ID, then the language-only reasoning stage is likely not at fault. From there, if a model exhibits sensitivity to masking the correct object region for \mathbf{K}_{wrong} but not for $\mathbf{K}_{correct}$, the vision encoder is the likely bottleneck. If there is no distinct sensitivity difference, the errors are likely taking place after the vision encoder, but before the linguistic reasoning. If model accuracy seems independent of both spatial ID correctness and image recognition capacity, the language model layers beyond spatial ID binding are likely the biggest bottleneck. Note that it is possible for incorrect spatial IDs to be correlated to wrong answers, but still have some model inaccuracies be resultant from factors other than spatial IDs, such as erroneous priors during LM readout (Leng et al., 2024; Ramakrishnan et al., 2018). In this case, it is still valuable to find if models can benefit from stronger spatial representations through this diagnosis process, and minimize avenues for failure. For the described analyses, we need a sufficient \mathbf{K}_{wrong} subset. As their \mathbf{K}_{wrong} are biggest on COCO-SPATIAL, we select LLaVA1.5 7B and LLaMA3.2VL 11B as model organisms for this section.

Ground Truth Spatial ID Deviation Experiment. First, we want to identify if models predict incorrect spatial IDs for the samples they get wrong. If the answer is *yes*, then it is likely that the downstream language model is not the performance bottleneck, since it is faithful to the spatial information received. To compute the deviation of the model’s believed spatial ID to the ground

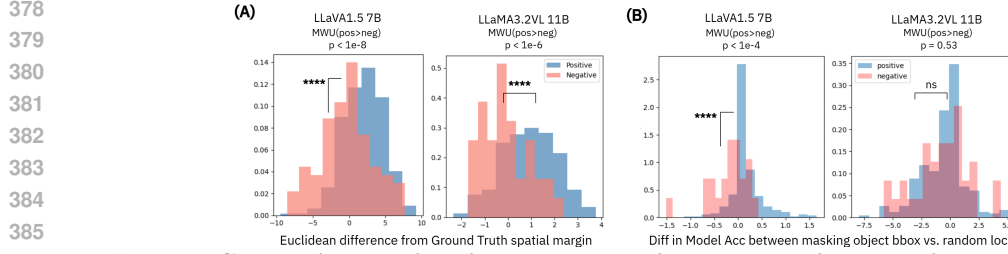


Figure 8: **Contrasting density histograms show incorrect spatial IDs drive bad predictions.** (A) shows deviation of model spatial IDs from g.t., and (B) the difference in model accuracy when masking objects vs. random locations in images. Histograms are samples VLMs got right (blue) or wrong (red). LLaVA shows faulty object detection with wrong answers, while LLaMA does not.

truth (g.t.), we compute the g.t. spatial ID by projecting the word activation onto the spatial axes:

$$\Delta_L^{(o)}(i, j)_{\text{ext}} \approx VV^T \phi_L(o; I_{(i,j)}^{(o)}, T^{(o)}), \quad V = [v_L, h_L] \quad (9)$$

For a spatial query like “Is the o to the left or right of a \bar{o} ?”, we can thus compute $\Delta^{(o)}(i, j)_{\text{gt}}$ and $\Delta^{(\bar{o})}(i, j)_{\text{gt}}$. The model’s assigned spatial IDs to the objects are computed per Eq. 9, for $\Delta^{(o)}(i, j)_{\text{ext}}$ and $\Delta^{(\bar{o})}(i, j)_{\text{ext}}$. Then the g.t. ID margin deviation for some object o is:

$$\text{ID deviation margin} = \epsilon_{\text{ext}} - \epsilon_{\text{gt}}, \quad \text{where } \epsilon_{\text{gt}} = i_{\text{gt}}^{(o)} - i_{\text{gt}}^{(\bar{o})}, \epsilon_{\text{ext}} = i_{\text{ext}}^{(o)} - i_{\text{ext}}^{(\bar{o})} \quad (10)$$

Here, a negative margin indicates that the model’s extracted spatial IDs oppose the ground truth.

ID Deviation Results. From Fig. 8A, we see that deviation from ground truth in extracted spatial ID margin is highly correlated with model mistakes. In other words, for LLaVA and LLaMA, wrong spatial IDs in object word activations led to wrong model answers, so linguistic reasoning was not the reason these failures occurred. Each subplot shows two density histograms overlaid in the same grid, where the x axis is $\epsilon_{\text{ext}} - \epsilon_{\text{gt}}$. The red histogram represents the density of ID deviations for $\mathbf{K}_{\text{wrong}}$, and the blue histogram shows the same for $\mathbf{K}_{\text{correct}}$. The red distribution is visibly skewed to the negatives compared to the blue. Quantitatively, we perform the Mann-Whitney U test (McKnight & Najab, 2010) to calculate the p-value for the hypothesis that the two distributions (red and blue) are non-identical. Now we ask, is this failure mode stemming from the vision encoder level, or does it occur during the spatial ID binding across modalities?

Image Masking Experiment. Altering the raw image input at the pixel level can help us understand whether it is a faulty vision encoder or faulty crossmodal information integration that has led to the failures. If the model’s beliefs on $\mathbf{K}_{\text{correct}}$ are more sensitive to masking the image raw input at the g.t. location of o , while beliefs on $\mathbf{K}_{\text{wrong}}$ change more with masking elsewhere, we can conclude that the vision encoder of this VLM is doing a poor job at object detection, leading to observed failures. If we do not observe this is the case, the failure may arise from the crossmodal information integration stage. In other words, the language model is doing a poor job appending binding IDs, despite the vision encoder having the necessary object recognition capacity.

We design an obfuscation paradigm inspired by methods like D-RISE (Petsiuk et al., 2021), where we either blur the bounding box of o , or R other locations in the image that do not intersect with the bboxes for o or \bar{o} . We then measure model belief change when masking the object vs. elsewhere:

$$\text{bbox sensitivity} = (P(\text{“GT”}) - P(\text{“GT”}|\text{mask } o)) - (P(\text{“GT”}) - \min_r [P(\text{“GT”}|\text{mask } r), r \in R]) \quad (11)$$

Image Masking Results. Fig. 8B shows overlaid histograms for bounding box masking sensitivities of $\mathbf{K}_{\text{correct}}$ and $\mathbf{K}_{\text{wrong}}$. Here, a negative value indicates greater sensitivity to raw pixel masking of random scenes, suggesting poor object detection. For LLaVA, there is a statistically significant p-value for the hypothesis that $\mathbf{K}_{\text{wrong}}$ is shifted more negative than $\mathbf{K}_{\text{correct}}$, indicating its vision encoder fails at object detection when it answers incorrectly. In contrast, $\mathbf{K}_{\text{wrong}}$, $\mathbf{K}_{\text{correct}}$ in LLaMA are agnostic to image obfuscation. This suggests that its failure modes likely stem after the vision encoder. These insights could be attributed to how LLaVA uses an out-of-the-box ViT that was text-aligned at a massive scale, hence not being tuned for finegrained detection, while LLaMA has a trained in-house ViT whose image-text alignment may be less robust.

Diagnosis Conclusion. With spatial IDs, we explore the causes for failure in a few model VLMs. We find that for both LLaMA and LLaVA, the linguistic reasoning stage is faithful to spatial IDs. LLaVA’s vision encoder is likely creating wrong spatial IDs from poor object detection, while LLaMA’s weak point appears to be information integration across the image patch activations to the text tokens. These conclusions are preliminary and do not suggest that *all* of a model’s failures stem from *one* architectural component, but can serve to guide finetuning stage choices when resources are scarce, or provide intuition for future model designs.

4.3 IMPROVING VLMs

Spatial IDs and Model Performance. To understand if spatial IDs could be a valuable learning signal, we first evaluate whether stronger steerability from spatial IDs is correlated to stronger models. Fig. 9 shows the results of this analysis, where indeed we see that models with higher zero-shot accuracy on COCO-spatial also exhibit greater belief changes with spatial ID interventions.

We define “steerability” as the difference between the change of belief resultant from steering with opposing spatial IDs versus noise. The layers of intervention are chosen as the middle third of all layers for that model. Each point shows the model’s mean steerability (on x) against its accuracy on COCO-spatial with no spatial intervention. Dotted lines connect models within a family.

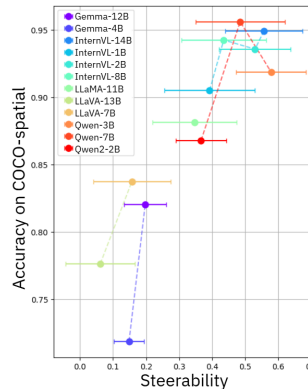


Figure 9: **Accuracy vs. Steerability.** Models with higher accuracy can be better steered with spatial IDs.

Spatial Loss Module

Fig. 9 shows spatial IDs signal stronger model performance. This suggests that the strength of spatial IDs could be a valuable learning signal for VLMs to learn principled spatial reasoning. To validate this intuition, we finetune Qwen2-2B on a synthetic dataset similar to the one used to extract spatial IDs, and evaluate on COCO-Spatial. We introduce an additional loss module at layer 11 that computes the cosine similarity between the predicted and ground-truth spatial ID at that layer. We provide detailed explanations for this process in §A.9. This spatial ID loss is added to the standard language modeling objective, providing extra supervision. We perform a control training without the spatial ID loss. Indeed, we see that explicit spatial ID loss helps the model generalize faster, reaching 90% accuracy on COCO-spatial at 3.2k steps, as shown below:

	Num Steps	0	800	1600	2400	3200
Control	LM Loss (↓)	3.30	0.05	<0.01	<0.01	<0.01
	COCO Val Accuracy (↑)	0.77	0.83	0.84	0.85	0.85
With Spatial Loss	LM Loss (↓)	2.71	0.04	<0.01	<0.01	<0.01
	Spatial ID Loss (↓)	0.75	0.58	0.41	0.36	0.33
	COCO Val Accuracy (↑)	0.77	0.83	0.84	0.88	0.91

5 TEMPORAL IDS IN VIDEO MODELS

Thus far, we have characterized spatial IDs as a causal model for spatial visual reasoning in VLMs. Could we find a similar linear paradigm for the temporal axis? In this section, we repeat the experiments in §2-3 for the temporal dimension in video models, with the goal of identifying linearly separable temporal markers on object words. For space, experimental procedures are described briefly here, and in greater detail in Appendix §A.

5.1 MIRRORING, EXTRACTING, AND STEERING ACROSS THE TEMPORAL AXIS

Temporal Mirror Swapping. We validate that there exist modality alignment layers with object-level visual information transfer in video models. For mirrored videos, we take the Scene_QA subset of MVBENCH (Li et al., 2024a) and swap the order of frames from back to front. Following Alg. 1, we show results for swapping text tokens, image patches, and object words in Fig. 10A. While

the error bound is noisier than spatial LLaVA, likely as LLaVA-Video follows response formats less well, we see the expected bump around middle layers for crossmodal integration.

Temporal ID Extraction. Derivation of temporal IDs and the temporal vector t_L follows Eq. 2 - 5, with synthetic 8-frame videos of OBJAVERSE renders. Results are shown in Fig. 10B. We again see that the text activation for “before” projects closer to earlier frames, than the activation for “after”.

Causality of Temporal IDs. Finally, to confirm controllability with arbitrary temporal IDs, we perform the steering experiment per Alg. 2 on MVBENCH videos. Results are shown in Fig. 10C. On these real, naturalistic videos, we see that later temporal IDs steer the model belief towards “after”, and earlier IDs towards “before”, as expected.

5.2 EMERGENCE OF TEMPORAL IDs

Fig. 10 shows summary results on LLaVA-Video, but we include temporal IDs from VideoLLaMA3 and Qwen2.5 in Appendix §B.2. LLaVA-Video and VideoLLaMA3 use textual description of the video length and number of frames to indicate timestamps preceding the visual input, while Qwen uses explicit MRoPE time IDs. This suggests that spatiotemporal IDs can emerge without explicit positional encoding, beyond the simple mechanism derived in Eq. 7.

6 RELATED WORK

Mechanistic interpretability is a growing field uncovering the inner workings of large models, popularizing techniques such as circuit tracing (Elhage et al., 2021; Ameisen et al., 2025), Sparse Autoencoders (Cunningham et al., 2023), linear probing (Alain & Bengio, 2016), and others. The Linear Representation Hypothesis posits that concepts are linearly encoded in LLM latents (Park et al., 2024), and activation patching supports that linear changes in activations drive model belief (Meng et al., 2022; Zhang & Nanda, 2023). Internal in-context reasoning mechanisms such as linear *binding IDs* (Feng & Steinhardt, 2024; Feng et al., 2024) have been identified in LLMs, along with other evidence for linear multi-hop reasoning (Yang et al., 2024), in-context task vectors (Hendel et al., 2023) and linear relational embeddings (Hernandez et al., 2024) during reasoning.

Linearity of embeddings have also been discovered in VLM latent spaces (Trager et al., 2023; Jiang et al., 2025a) to some degree. Previous work showed that VLMs separate VQA into image-focused then text-focused stages (Jiang et al., 2025b), and others have extended LLM interpretability techniques like logit lens (Neo et al., 2024) or attention tracking (Zhang et al., 2024a; Yu & Ananiadou, 2025) to VLMs to unearth internal circuits. In our work, we mechanistically capture spatiotemporal information flow from image patches to text tokens in VLMs, via the spatial ID mechanism.

7 CONCLUSION, LIMITATIONS, & FUTURE WORK

We propose spatiotemporal IDs as a linear model for visual reasoning about space and time in VLMs. With a series of causal analyses, we show these IDs can be obtained in many SoTA models, and that they closely mediate models’ beliefs about visual objects’ location in space and time. We further offer ways to extend this mechanistic insight to improving existing VLMs. For tractability, our work is currently limited to analyses in simple spatial queries or appearance-based temporal queries. Experimental design for more complex, open-ended queries will enhance our understanding of how VLMs utilize rudimentary concepts like spatial IDs in more nuanced settings. Further, we only extract and steer models of sizes up to 14B parameters due to compute constraints. Investigation into whether the spatial ID circuit plays a similarly prominent role in larger models will reveal whether VLMs of varying capacities follow analogous methods for visual reasoning, or employ distinct measures. Lastly, while we show two potential ways to leverage spatial IDs for VLM diagnostics, future work could include other use cases, such as spatiotemporal IDs as a proxy learning signal, a motivator for explicit temporal encodings, and more.

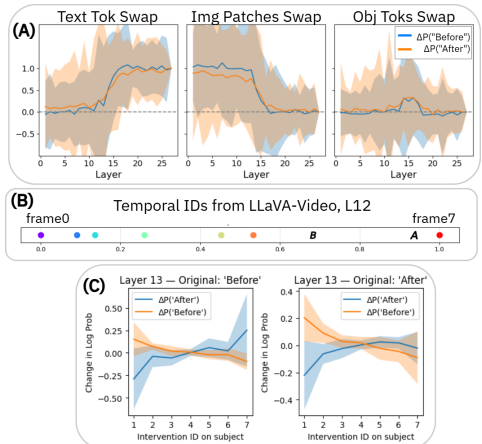


Figure 10: **Temporal ID Results.** Mirror Swapping on videos (A), Temporal ID grid (B), and temporal ID steering on model beliefs (C) with LLaVA-Video

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

REPRODUCIBILITY

We provide finegrained details for all experiments in §A of the Appendix, and results on all the models considered in §B. We include experimental details, results from various ablation analyses and counterfactual trials in §C-D. We will release the code for reproducing all results upon acceptance.

REFERENCES

- Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes. *arXiv preprint arXiv:1610.01644*, 2016.
- Emmanuel Ameisen, Jack Lindsey, Adam Pearce, Wes Gurnee, Nicholas L. Turner, Brian Chen, Craig Citro, David Abrahams, Shan Carter, Basil Hosmer, Jonathan Marcus, Michael Sklar, Adly Templeton, Trenton Bricken, Callum McDougall, Hoagy Cunningham, Thomas Henighan, Adam Jermy, Andy Jones, Andrew Persic, Zhenyi Qi, T. Ben Thompson, Sam Zimmerman, Kelley Rivoire, Thomas Conerly, Chris Olah, and Joshua Batson. Circuit tracing: Revealing computational graphs in language models. *Transformer Circuits Thread*, 2025. URL <https://transformer-circuits.pub/2025/attribution-graphs/methods.html>.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-VL Technical Report, February 2025. URL <http://arxiv.org/abs/2502.13923>. arXiv:2502.13923 [cs].
- Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14455–14465, 2024a.
- Shiqi Chen, Tongyao Zhu, Ruochen Zhou, Jinghan Zhang, Siyang Gao, Juan Carlos Niebles, Mor Geva, Junxian He, Jiajun Wu, and Manling Li. Why is spatial reasoning hard for vlms? an attention mechanism perspective on focus areas. *arXiv preprint arXiv:2503.01773*, 2025.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. InternVL: Scaling up Vision Foundation Models and Aligning for Generic Visual-Linguistic Tasks, January 2024b. URL <http://arxiv.org/abs/2312.14238>. arXiv:2312.14238 [cs].
- Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*, 2023.
- Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 13142–13153, 2023.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pp. arXiv–2407, 2024.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. <https://transformer-circuits.pub/2021/framework/index.html>.
- Yue Fan, Xuehai He, Diji Yang, Kaizhi Zheng, Ching-Chen Kuo, Yuting Zheng, Sravana Jyothi Narayanaraju, Xinze Guan, and Xin Eric Wang. GRIT: Teaching MLLMs to Think with Images, May 2025. URL <http://arxiv.org/abs/2505.15879>. arXiv:2505.15879 [cs].

- 594 Jiahai Feng and Jacob Steinhardt. How do Language Models Bind Entities in Context?, May 2024.
595 URL <http://arxiv.org/abs/2310.17191>. arXiv:2310.17191 [cs].
596
- 597 Jiahai Feng, Stuart Russell, and Jacob Steinhardt. Monitoring Latent World States in Language
598 Models with Propositional Probes, December 2024. URL [http://arxiv.org/abs/2406.](http://arxiv.org/abs/2406.19501)
599 19501. arXiv:2406.19501 [cs].
- 600 Roei Hendel, Mor Geva, and Amir Globerson. In-Context Learning Creates Task Vectors, October
601 2023. URL <http://arxiv.org/abs/2310.15916>. arXiv:2310.15916 [cs].
602
- 603 Evan Hernandez, Arnab Sen Sharma, Tal Haklay, Kevin Meng, Martin Wattenberg, Jacob Andreas,
604 Yonatan Belinkov, and David Bau. Linearity of Relation Decoding in Transformer Language
605 Models, February 2024. URL <http://arxiv.org/abs/2308.09124>. arXiv:2308.09124
606 [cs].
- 607 Nick Jiang, Anish Kachinthaya, Suzie Petryk, and Yossi Gandelsman. Interpreting and Editing
608 Vision-Language Representations to Mitigate Hallucinations, February 2025a. URL [http://](http://arxiv.org/abs/2410.02762)
609 arxiv.org/abs/2410.02762. arXiv:2410.02762 [cs].
610
- 611 Zhangqi Jiang, Junkai Chen, Beier Zhu, Tingjin Luo, Yankun Shen, and Xu Yang. Devils in Middle
612 Layers of Large Vision-Language Models: Interpreting, Detecting and Mitigating Object Hallu-
613 cinations via Attention Lens, April 2025b. URL <http://arxiv.org/abs/2411.16724>.
614 arXiv:2411.16724 [cs].
- 615 Amita Kamath, Jack Hessel, and Kai-Wei Chang. What’s” up” with vision-language models? in-
616 vestigating their struggle with spatial reasoning. *arXiv preprint arXiv:2310.19785*, 2023.
617
- 618 Raphi Kang, Yue Song, Georgia Gkioxari, and Pietro Perona. Is clip ideal? no. can we fix it? yes!
619 *arXiv preprint arXiv:2503.08723*, 2025.
- 620 Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing.
621 Mitigating object hallucinations in large vision-language models through visual contrastive de-
622 coding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recogni-*
623 *tion*, pp. 13872–13882, 2024.
624
- 625 Kunchang Li, Yali Wang, Yanan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo
626 Chen, Ping Lou, Limin Wang, and Yu Qiao. MVBench: A Comprehensive Multi-modal
627 Video Understanding Benchmark. In *2024 IEEE/CVF Conference on Computer Vision and*
628 *Pattern Recognition (CVPR)*, pp. 22195–22206, Seattle, WA, USA, June 2024a. IEEE. ISBN
629 9798350353006. doi: 10.1109/CVPR52733.2024.02095. URL [https://ieeexplore.](https://ieeexplore.ieee.org/document/10658165/)
630 [ieee.org/document/10658165/](https://ieeexplore.ieee.org/document/10658165/).
- 631 Lei Li, Yuanxin Liu, Linli Yao, Peiyuan Zhang, Chenxin An, Lean Wang, Xu Sun, Lingpeng Kong,
632 and Qi Liu. Temporal reasoning transfer from text to video. *arXiv preprint arXiv:2410.06166*,
633 2024b.
- 634
- 635 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
636 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European*
637 *conference on computer vision*, pp. 740–755. Springer, 2014.
- 638 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual Instruction Tuning, December
639 2023. URL <http://arxiv.org/abs/2304.08485>. arXiv:2304.08485 [cs].
640
- 641 Patrick E McKnight and Julius Najab. Mann-whitney u test. *The Corsini encyclopedia of psychology*,
642 pp. 1–1, 2010.
- 643 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual
644 associations in gpt. *Advances in neural information processing systems*, 35:17359–17372, 2022.
645
- 646 Clement Neo, Luke Ong, Philip Torr, Mor Geva, David Krueger, and Fazl Barez. Towards In-
647 terpreting Visual Information Processing in Vision-Language Models, October 2024. URL
<http://arxiv.org/abs/2410.07149>. arXiv:2410.07149 [cs].

- 648 Kiho Park, Yo Joong Choe, and Victor Veitch. The Linear Representation Hypothesis and the Geom-
649 etry of Large Language Models, July 2024. URL <http://arxiv.org/abs/2311.03658>.
650 arXiv:2311.03658 [cs].
- 651 Vitali Petsiuk, Rajiv Jain, Varun Manjunatha, Vlad I Morariu, Ashutosh Mehra, Vicente Ordonez,
652 and Kate Saenko. Black-box explanation of object detectors via saliency maps. In *Proceedings of*
653 *the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11443–11452, 2021.
- 654 Sainandan Ramakrishnan, Aishwarya Agrawal, and Stefan Lee. Overcoming language priors in
655 visual question answering with adversarial regularization. *Advances in neural information pro-*
656 *cessing systems*, 31, 2018.
- 657 Bin Ren, Yahui Liu, Yue Song, Wei Bi, Rita Cucchiara, Nicu Sebe, and Wei Wang. Masked jigsaw
658 puzzle: A versatile position embedding for vision transformers. In *Proceedings of the IEEE/CVF*
659 *Conference on Computer Vision and Pattern Recognition*, pp. 20382–20391, 2023.
- 660 Ilias Stogiannidis, Steven McDonagh, and Sotirios A Tsaftaris. Mind the gap: Benchmarking spatial
661 reasoning in vision-language models. *arXiv preprint arXiv:2503.19707*, 2025.
- 662 Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya
663 Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open
664 models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- 665 Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide
666 shut? exploring the visual shortcomings of multimodal llms. In *Proceedings of the IEEE/CVF*
667 *Conference on Computer Vision and Pattern Recognition*, pp. 9568–9578, 2024.
- 668 Matthew Trager, Pramuditha Perera, Luca Zancato, Alessandro Achille, Parminder Bhatia, and Ste-
669 fano Soatto. Linear spaces of meanings: compositional structures in vision-language models. In
670 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 15395–15404,
671 2023.
- 672 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
673 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information*
674 *Processing Systems (NeurIPS)*, 2017. URL <https://arxiv.org/abs/1706.03762>.
- 675 Junbin Xiao, Angela Yao, Yicong Li, and Tat-Seng Chua. Can i trust your answer? visually grounded
676 video question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
677 *Pattern Recognition*, pp. 13204–13214, 2024.
- 681 Sohee Yang, Elena Gribovskaya, Nora Kassner, Mor Geva, and Sebastian Riedel. Do Large Lan-
682 guage Models Latently Perform Multi-Hop Reasoning?, February 2024. URL <http://arxiv.org/abs/2402.16837>.
683 arXiv:2402.16837 [cs].
- 684 Zeping Yu and Sophia Ananiadou. Understanding Multimodal LLMs: the Mechanistic Interpretabil-
685 ity of Llava in Visual Question Answering, January 2025. URL <http://arxiv.org/abs/2411.10950>.
686 arXiv:2411.10950 [cs].
- 687 Boqiang Zhang, Kehan Li, Zesen Cheng, Zhiqiang Hu, Yuqian Yuan, Guanzheng Chen, Sicong
688 Leng, Yuming Jiang, Hang Zhang, Xin Li, Peng Jin, Wenqi Zhang, Fan Wang, Lidong Bing, and
689 Deli Zhao. VideoLLaMA 3: Frontier Multimodal Foundation Models for Image and Video Under-
690 standing, January 2025. URL <http://arxiv.org/abs/2501.13106>. arXiv:2501.13106
691 [cs].
- 692 Fred Zhang and Neel Nanda. Towards best practices of activation patching in language models:
693 Metrics and methods. *arXiv preprint arXiv:2309.16042*, 2023.
- 694 Xiaofeng Zhang, Yihao Quan, Chen Shen, Xiaosong Yuan, Shaotian Yan, Liang Xie, Wenxiao
695 Wang, Chaochen Gu, Hao Tang, and Jieping Ye. From Redundancy to Relevance: Informa-
696 tion Flow in LVLMS Across Reasoning Tasks, October 2024a. URL <http://arxiv.org/abs/2406.06579>.
697 arXiv:2406.06579 [cs].
- 698 Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video
699 Instruction Tuning With Synthetic Data, October 2024b. URL <http://arxiv.org/abs/2410.02713>.
700 arXiv:2410.02713 [cs].

702	LINEAR MECHANISMS FOR SPATIOTEMPORAL REASONING IN VISION	
703	LANGUAGE MODELS	
704	–SUPPLEMENTARY MATERIAL–	
705		
706		
707		
708	A Experimental Details	15
709	A.1 Mirror Swapping	15
710	A.2 Spatiotemporal ID Extraction	16
711	A.3 Arbitrary Steering Experiments	17
712	A.4 Color-binding Reasoning Experiments	17
713	A.5 Adversarial Steering Experiments	18
714	A.6 ID Deviation	19
715	A.7 Obfuscating Experiments	19
716	A.8 Model Diagnosis Addendum	19
717	A.9 Model finetuning with Spatial Loss	20
718		
719		
720		
721		
722		
723		
724		
725	B Experimental Results on More Models	21
726	B.1 Spatial Grids on More Models	21
727	B.2 Temporal Grids on more Models	23
728		
729		
730		
731		
732	C Counterfactuals	23
733	C.1 Spatial IDs from non-object words	23
734	C.2 Mirror Swapping on non-object words	24
735	C.3 Steering Effects on Orthogonal Directions (x vs. y), (time vs. x)	25
736		
737		
738		
739	D Ablations	26
740	D.1 Scaling Analysis for Spatiotemporal ID Extraction	26
741	D.2 Varying prompt wording and object sizes during extraction	26
742		
743		
744		
745		
746	E Theoretical Analysis of Spatial IDs	27
747	E.1 Informal Proof for Spatial ID Emergence	27
748	E.2 Empirical Relationship between Positional Encoding and Spatial IDs	28
749		
750		
751		
752	F LLM Usage Disclosure	29
753		
754		
755		

A EXPERIMENTAL DETAILS

A.1 MIRROR SWAPPING

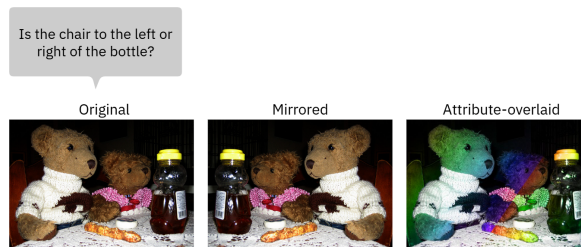


Figure A1: Example altered images for mirror swapping and attribute swapping.

Token handling. Different models and tokenizers have different tokenizing schemes. For example, for the query “Question: Is the the thermometer to the left or right of the desktop? Answer left or right. Answer: ”, the tokenization from Gemma, LLaMA, LLaVA, and Qwen will be as shown:

=== Gemma ===

```
Tokens: ['Question', ':', '_Is', '_the', '_thermometer', '_to',
'_the', '_left', '_or', '_right', '_of', '_the', '_desktop', '?',
'_Answer', '_left', '_or', '_right', '.', '_Answer', ':']
```

=== LLaMA ===

```
Tokens: ['_Question', ':', '_Is', '_the', '_therm', 'ometer',
'_to', '_the', '_left', '_or', '_right', '_of', '_the', '_desktop',
'?', '_Answer', '_left', '_or', '_right', '.', '_Answer', ':']
```

=== LLaVA ===

```
Tokens: ['_Question', ':', '_Is', '_the', '_therm', 'ometer',
'_to', '_the', '_left', '_or', '_right', '_of', '_the', '_desktop',
'?', '_Answer', '_left', '_or', '_right', '.', '_Answer', ':']
```

=== Qwen ===

```
Tokens: ['Question', ':', 'GIs', 'Gthe', 'Gthermometer', 'Gto',
'Gthe', 'Gleft', 'Gor', 'Gright', 'Gof', 'Gthe', 'Gdesktop', '?',
'GAnswer', 'Gleft', 'Gor', 'Gright', '.', 'GAnswer', ':']
```

When a word is represented as multiple tokens per a model’s processor (e.g., *frog* is tokenized into $[_f, rog]$ in LLaVA), we take the last index of this list to be most representative of the object, as it is the distinguishing element. So in the case of LLaMA or LLaVA, we would take the “ometer” token to represent the object “thermometer”.

Logit Probabilities. For assessing the model’s likelihood of saying “left” vs. “right”, or two other options, we take the log probability for that token following the tokenization scheme of the model family. This means we take the model output.logits and index at the token ID for ‘Gright’ in Qwen, for example, to get $P(\text{“right”})$.

Activation Patching. For every model, we first register a forward hook at each layer to collect the intermediate activation for both the original (case 1) and mirror-swapped (case 2) cases. Then, we register another forward-hook replace the original activations with the mirror-swapped one at select indices according to the three different settings.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

A.2 SPATIOTEMPORAL ID EXTRACTION

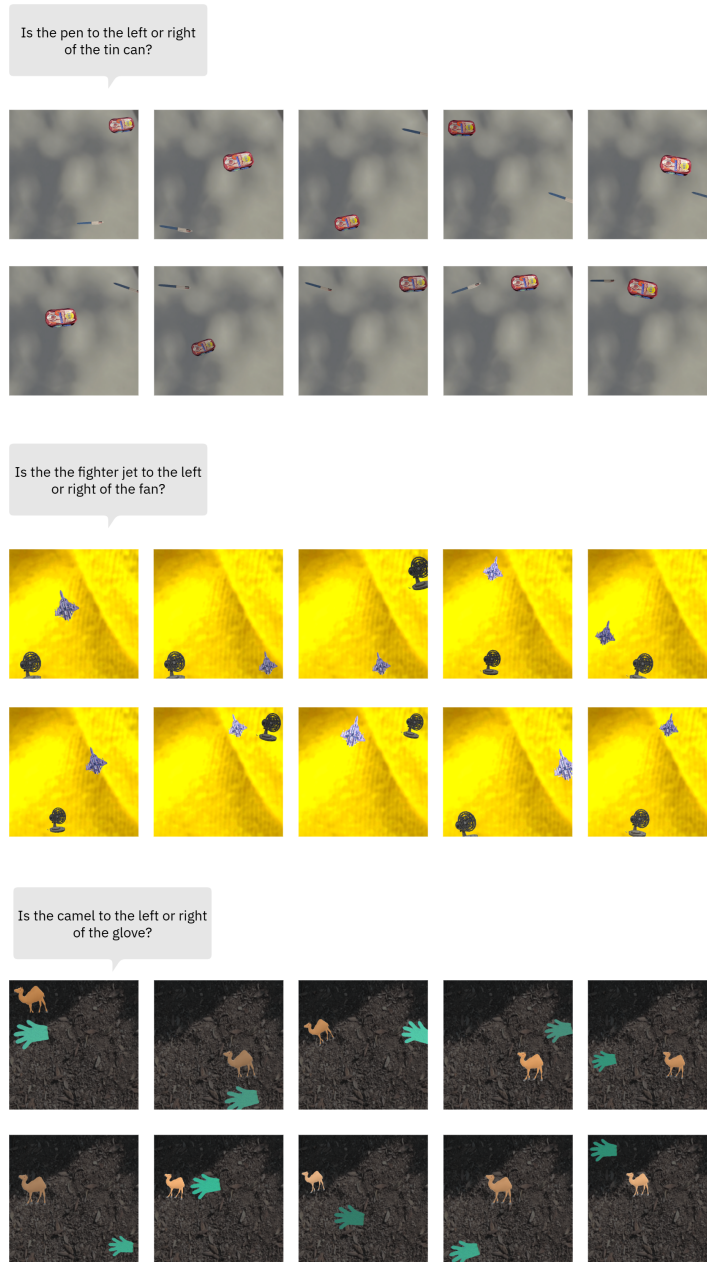


Figure A2: Synthetic images used towards spatial ID extraction

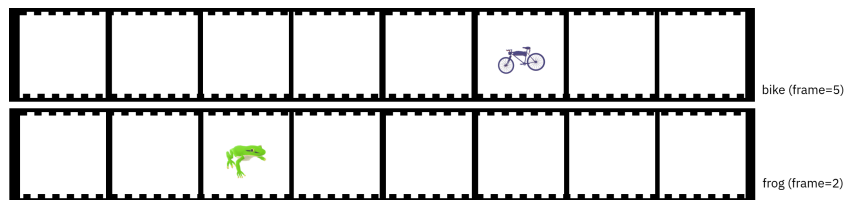


Figure A3: Illustration of synthetic videos used towards temporal ID generation. All videos had 8 frames.

Synthetic Image generation for Spatial IDs. We use 55 OBJAVERSE object renderers and project them in various pairs onto random backgrounds, per (Kang et al., 2025). All images with the same objects get the same text query. In §D.1 we show the different number of objects and total number of images used to generate spatial IDs. For w object pairs, we generate $w \times s \times m^2 \times (m^2 - 1)$ total images, where s is the number of object sizes we consider, and m is grid size. While we find minimal difference with extraction dataset size, as shown in §D.1, we use 90 object pairs, and consider $s = 4$ from $\{224, 174, 124, 74\}$, yielding 86,400 images. Note that each image size is $224 \times 4 = 896$ in width and height.

Synthetic Video generation for Temporal IDs. We take 5 unique OBJAVERSE object pairs in 61 distinct temporal arrangements. For baseline temporal ID extraction, all objects were centered in the image. For spatial vs. temporal disentanglement verification, we try three spatial variants - left, center, and right - for object location.

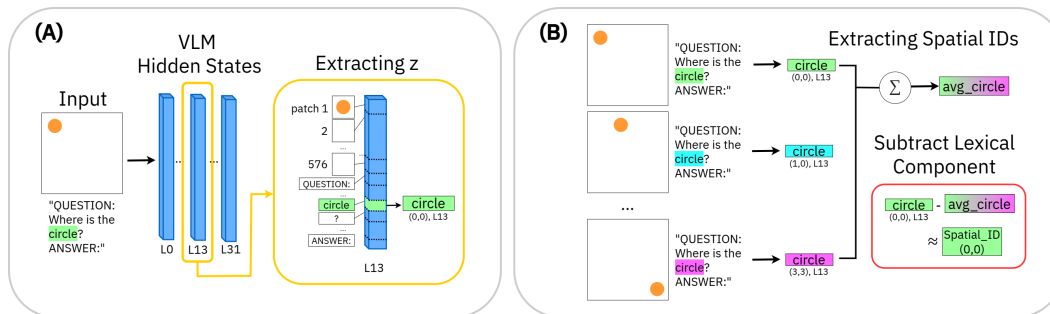


Figure A4: Illustration of spatial ID extraction. We isolate the relevant visual object word token in a chosen layer activation (A) and compute the shared lexical component for that particular object word that is independent of spatial localization (B) to acquire the linearly bound spatial ID.

A.3 ARBITRARY STEERING EXPERIMENTS

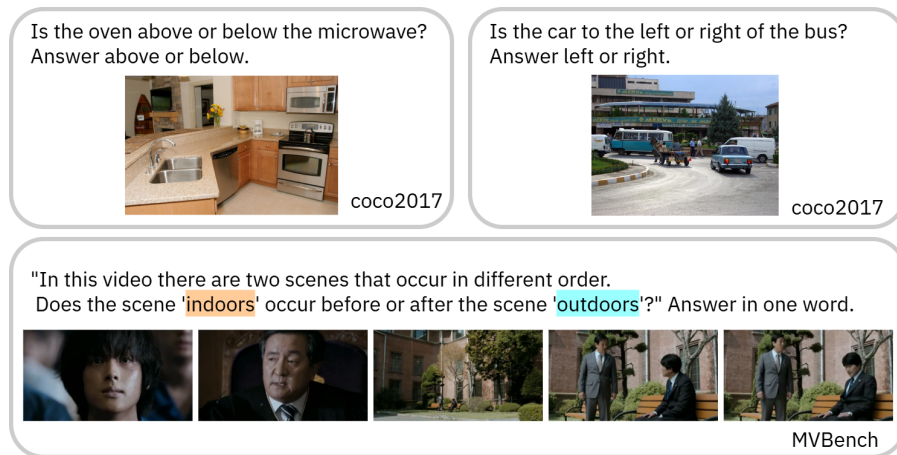


Figure A5: Examples of real images (top) and videos (bottom) we use to test model beliefs.

A.4 COLOR-BINDING REASONING EXPERIMENTS

Do spatial IDs mediate visual reasoning beyond direct spatial queries (such as A above/below B, etc.)? To test this, we perform mirror swapping on two images where two objects are *in the same place*, unlike the mirror swapping in §2.1. This time, the objects are opposing in color. Fig. A6 shows the example query setup, as well as the results of swapping all image tokens, all text tokens, just the color word tokens, or just the object word token.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935

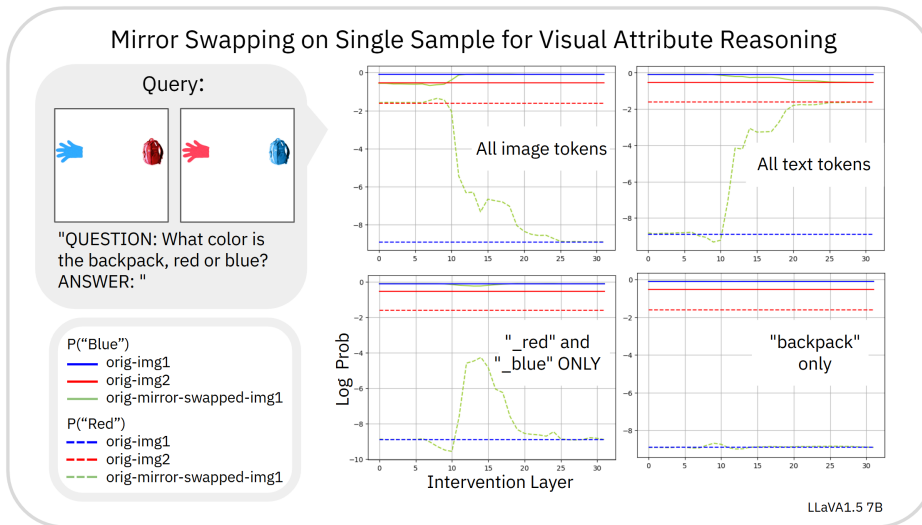


Figure A6: Mirror Swapping on Single Sample for Color Binding

936
937
938
939

Notice that swapping the activations for “backpack” has no effect, since the spatial ID encoded in the object word activation stays the same regardless of the input image (the backpack is in the same location in both images). Swapping the activations of color-related words, on the other hand, alters model belief at key modality binding layers. This suggests that the color words were storing spatial IDs that corresponded to the location where that color was present, and matching this color spatial ID to the object token was the readout process.

945
946
947
948

We repeat this experiment across 100 total such images, and show the results in Fig. A7. On average we see that swapping the color word tokens influences model beliefs in intermediate layers, much more so than swapping non-color word tokens. This suggests that spatial IDs mediate visual reasoning beyond direct spatial queries.

949
950
951
952
953
954
955
956
957
958

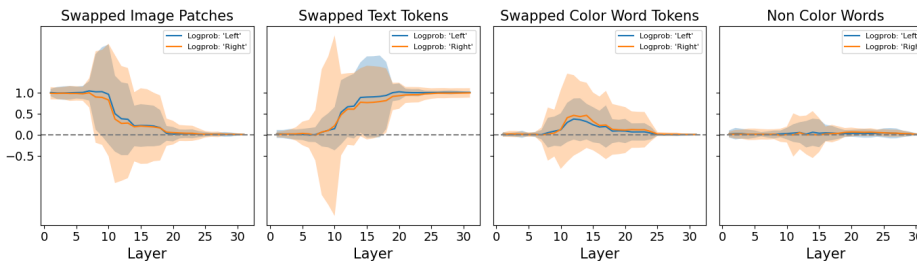


Figure A7: Swapping tokens for color binding.

959
960
961
962
963

A.5 ADVERSARIAL STEERING EXPERIMENTS

965
966
967
968
969
970
971

We perform steering on layers 9 through layer $2(\text{model.len} // 3)$ per model. To ensure activation norms don’t explode, we test a few different scaling factors for the intervening spatial ID’s norm. In the scaling factor = 1 case, we scale the norm of the spatial ID to equal the norm of that word token’s activation vector. We try a few scaling factors and choose 5 for steering all models, both for the noise vectors as well as the spatial IDs. Here, the norm of the spatial ID is fine to exceed that of the original token activation, as we subtract the opposing spatial ID to readjust the norm. This is shown in Alg. 2. For confidence intervals, we choose the three layers which had greatest steering effect, and report equivalent layers’ effects for the noise case.

A.6 ID DEVIATION

Classifying Model Belief. For the ID deviation experiment, we classified the model’s belief based on its decoded response. If the response contains only the correct spatial relationship (e.g. left) and not the incorrect spatial relationship, it’s considered correct. If the response contains only the incorrect spatial relationship, it’s considered incorrect. If none or both were present, it’s considered nonsense and discarded.

A.7 OBFUSCATING EXPERIMENTS

We take images from COCO_SPATIAL and Gaussian blur different regions, as below. We generate 4 images blurred in incorrect regions in addition to the 1 image with the correct bounding box blurred. For the sensitivity, we take the difference between the outside region which changed the model belief the most, and the bounding box.



Figure A8: Example of original query and two blurred options. Yellow grid lines are just for visualization.

A.8 MODEL DIAGNOSIS ADDENDUM

Oracle Injection Experiment To further isolate what model components may be responsible for creating incorrect spatial IDs, we conduct the oracle injection experiment. Specifically, we intervene with the *correct* spatial IDs on the object words at different layers, and see how that changes model accuracy from the control case without any intervention.

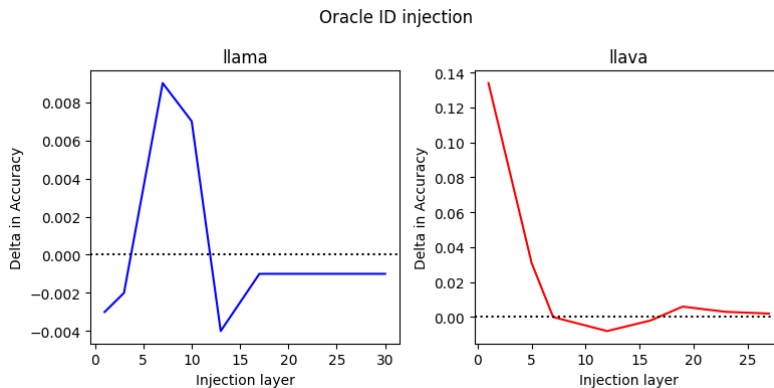


Figure A9: LLaMA and LLaVA evaluation accuracy on synthetic grid-like data with oracle spatial ID injections at varying layers. 0 is baseline model performance, without any intervention.

In accordance with our preliminary conclusion from §4.2, we see that LLaVA models’ accuracy increases 13.4% above the baseline when injected with oracle truth spatial IDs at layer 1. This suggests that indeed, if the image encoder had supplied correct spatial information, the downstream LM of LLaVA would have yielded greater accuracy. Intervention on intermediate to later layers in LLaVA has little effect. In LLaMA, we see that intervening on the earliest layers actually has little effect, while intervening on intermediate layers preceding the modality integration layer increases model accuracy by a modest amount (1%). Note that the low percentage is likely because LLaMA

has higher accuracy on this spatial dataset to begin with. This behavior is in line with our expectation from §4.2, where we do not expect it to benefit greatly from altering image encoder spatial localization performance, but instead benefit from spatial information condensation into the proper tokens.

For this experiment, we evaluated on synthetic images made with objaverse, such as those shown in Fig. A2. The interventions were performed with IDs from layer 17 on LLaMA for all layers including and below 17, and IDs from layer 12 on LLaVA for all layers including and below 12, as these were the layers identified as carrying spatial ID information in these respective models. LLaMA interventions were performed at layers [1, 3, 7, 10, 13, 17, 21, 25, 30, 35] and LLaVA interventions on [1, 5, 7, 12, 16, 19, 23, 27].

A.9 MODEL FINETUNING WITH SPATIAL LOSS

Spatial ID Loss Module In §4.3 we described finetuning Qwen2-2B with a spatial ID augmented loss module. Specifically, we freeze all weights except the MLPs of the last six vision encoder blocks, which we believe are most important for spatial reasoning, and train with synthetic data akin to those shown in Fig. A2. We batch 15 images of the same object but varying locations into a mini-batch, and compute the predicted spatial ID by subtracting the average activation. This is similar to how we extracted the spatial IDs in §2.2.

The validation accuracy and training plots are shown in Fig. A10. We see that with spatial ID loss, model accuracy on the naturalistic validation set (COCO-spatial) increases around 6% (absolute) beyond the baseline plateau, reaching a 90% accuracy in under 2.8k steps.

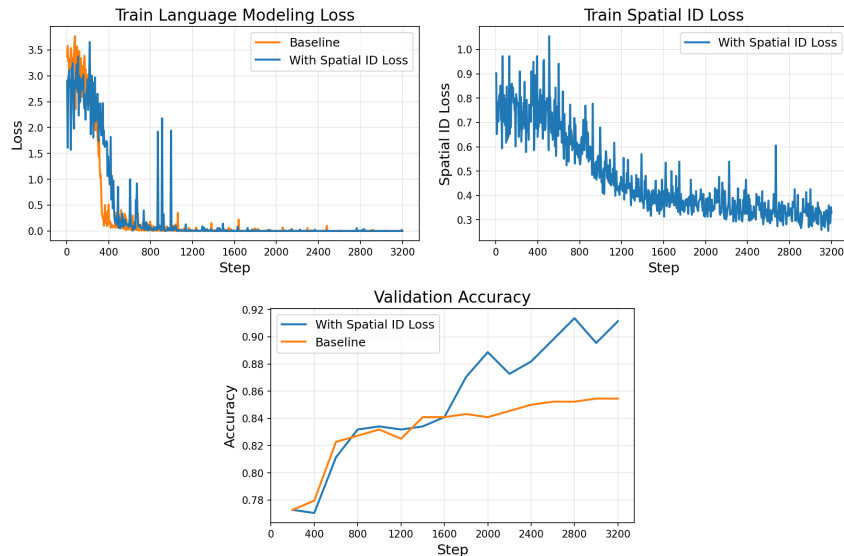


Figure A10: Plots from Qwen2-2B finetuning with and without spatial ID loss

B EXPERIMENTAL RESULTS ON MORE MODELS

B.1 SPATIAL GRIDS ON MORE MODELS

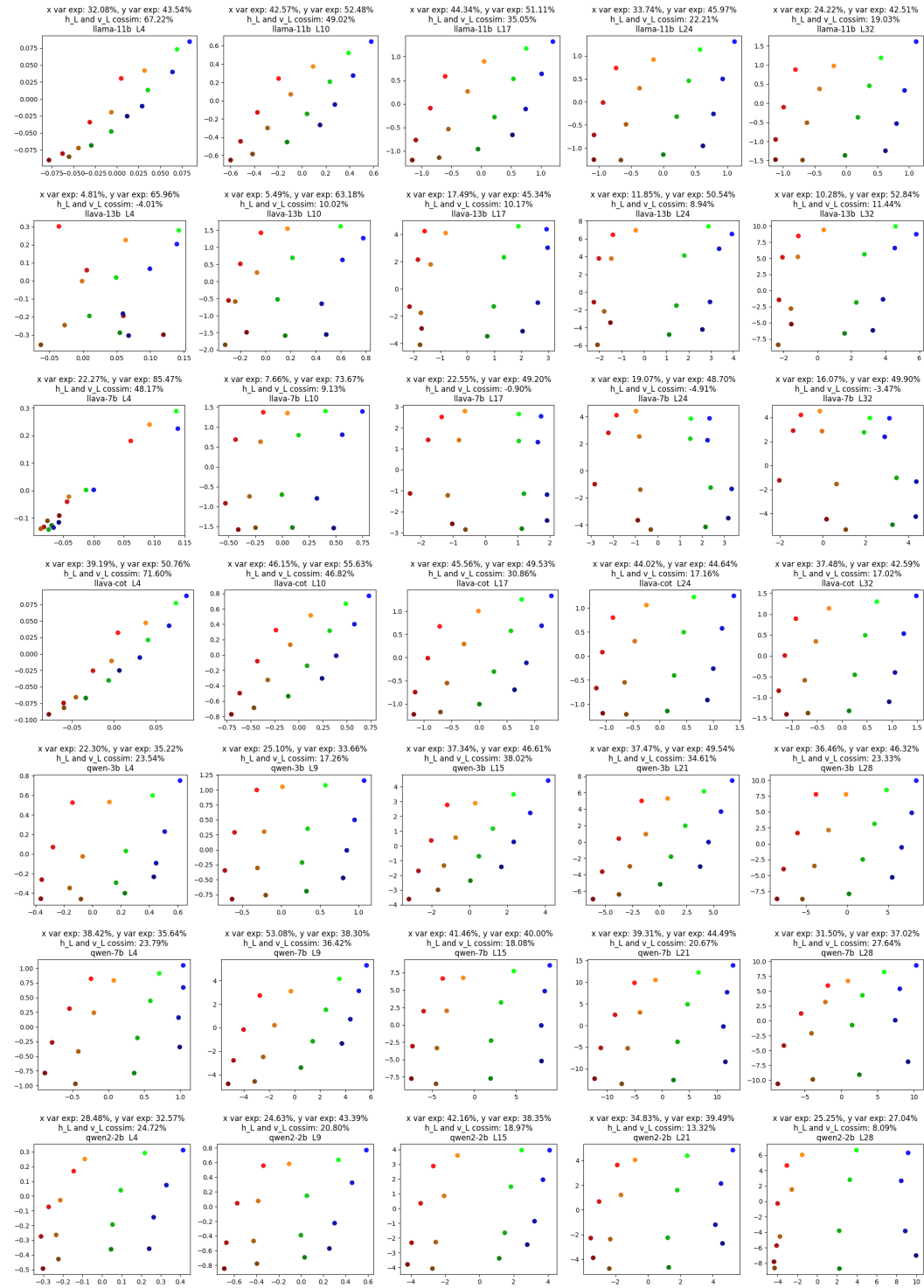


Figure A11: Spatial ID grids for LLaVA, LLaMA, and Qwen models.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

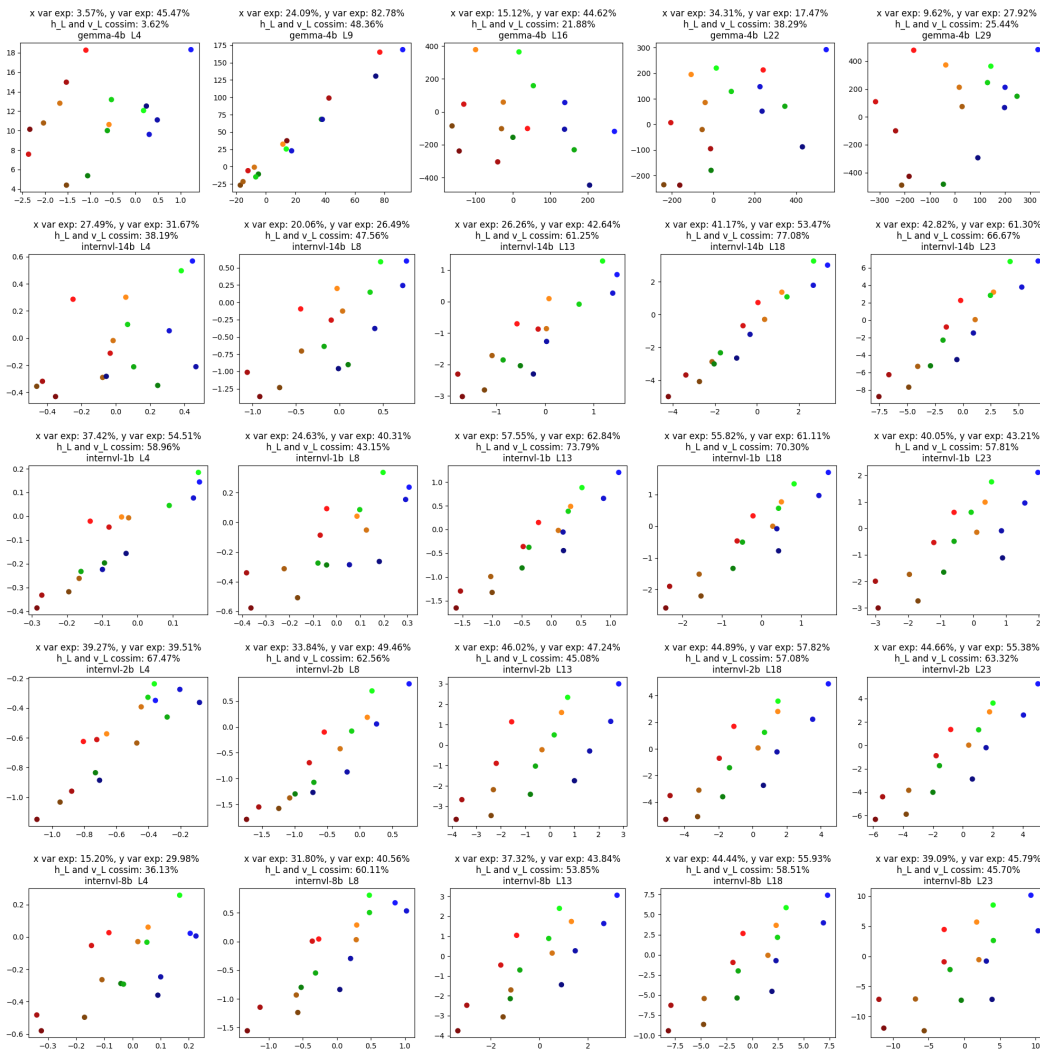


Figure A12: Spatial ID grids for Gemma and InternVL models.

Fig. A11,A12 show spatial ID grids for all models shown in Fig. 2. Subplot headings include % variance explained by each spatial axis, as well as the cosine similarity between the spatial axes. Notably, spatial IDs on some models seem to yield highly correlated v_L and h_L , suggesting different spatial directions may be conflated.

B.2 TEMPORAL GRIDS ON MORE MODELS

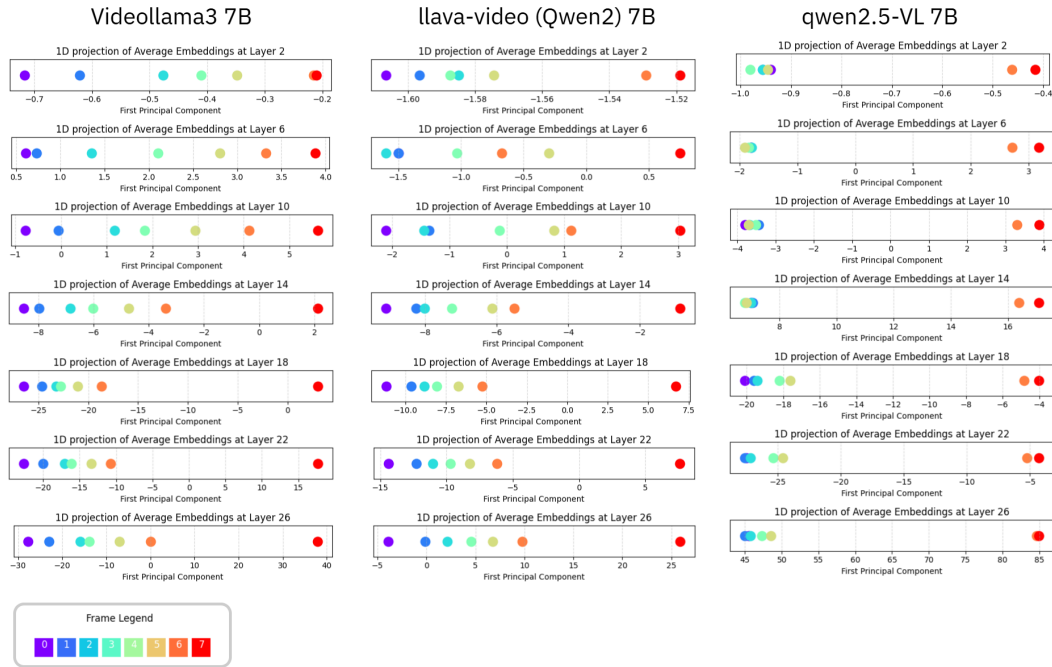


Figure A13: Temporal ID grids for VideoLLaMA3, LLaVA-Video, and Qwen2.5.

Across the models, there is a trend for the last frame(s) to be much farther away from the rest of the frames’ temporal IDs. This may be a result of the data bias during model training, where a lot of instruction tuning datasets will ask temporal questions that only require paying attention to the last frame (e.g., *did the person leave the room?* only requires looking at the first and last frame, and intermediate nuances are less important).

C COUNTERFACTUALS

C.1 SPATIAL IDS FROM NON-OBJECT WORDS

In §2.1 we concluded that spatial information is largely stored in object words at intermediate layers. But could the information storage be spread out across the sequence dimension in internal activations? To test this, we extract spatial IDs from non-object words, per §2.2. Specifically we choose the spatial words in the query format “Is the {obj_word1} {spatial_word1} {spatial_word2} {obj_word2}?”.

We then perform steering on object words, as well as non-object words, with both the spatial IDs extracted from object words and non-object words. We use the same steering algorithm as Alg. 2. The results are shown in Table 1. We see that some spatial semantics seems to be extractable from non-object words, and model belief is partially steerable through non-object words when using spatial IDs from object words, likely due to the fact that semantic word meanings are rarely perfectly contained within the initial word token in practical applications. In particular, spatial word tokens are likely to have information bleed over from the object word tokens while performing spatial queries, due to the way attention merges information between similar sequences. Regardless, effects from steering on object words with spatial IDs from object words is by far the strongest.

Model Name	ID-LR/apply-LR	ID-LR / apply-obj	ID-obj / apply-obj	ID-obj / apply-LR
Qwen-3b	18.77	51.19	81.23	48.46
Qwen-7b	6.12	38.78	72.35	47.96
LLaVA-7b	29.83	30.51	46.26	26.19
LLaVA-13b	8.16	19.39	48.30	15.25

Table 1: Spatial IDs extracted from object words and applied onto object words are most successful at steering model beliefs. Spurious effects are observed from IDs extracted from or applied unto unrelated words, but the effects are clearly concentrated on the object words.

C.2 MIRROR SWAPPING ON NON-OBJECT WORDS

To first showcase on a single sample the difference between mirror swapping on object tokens versus non object tokens, we choose a synthetic example with two objects on a blank background. Fig. A14 shows the results. Here, the green line shows model belief change, and the x axis indicates the layer of intervention. Mirror swapping at just the object words has a slightly less prominent effect than intervening at the object words in addition to immediate neighboring tokens (such as the space preceding the animal word), which captures some of the information bleed. Swapping all tokens except for those belonging to object words, on the other hand, has the smallest observable effect. Hence spatial information is likely concentrated in object tokens.

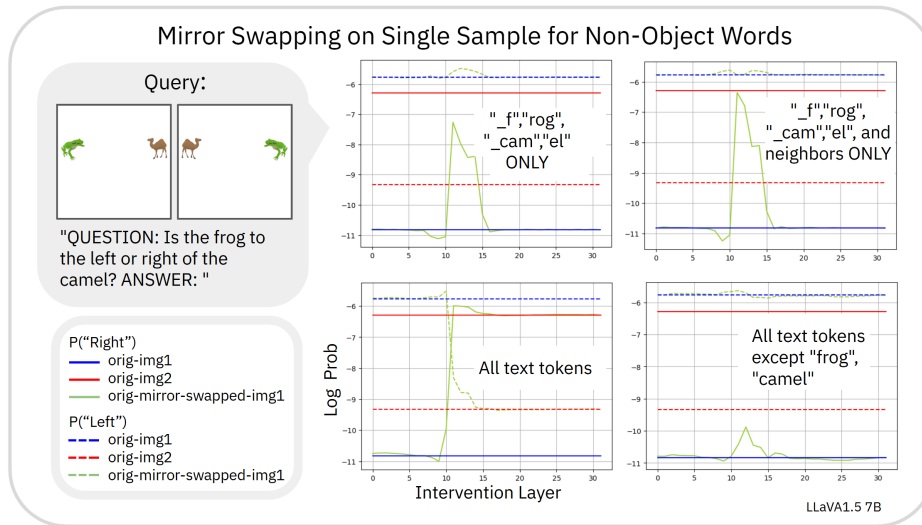


Figure A14: Synthetic image example for mirror swapping. Swapping non-object tokens has minimal impact on model belief change.

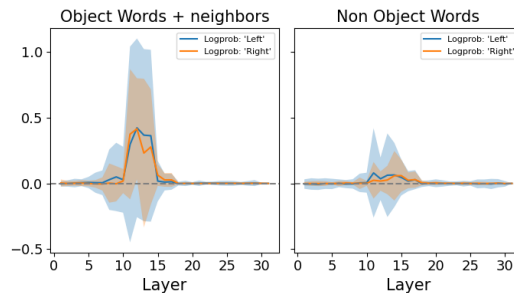


Figure A15: Mirror Swapping non object words

We can now repeat the mirror swapping at non-object tokens at scale on COCO images. Fig. A15 shows the difference between steering on object words and immediate neighboring tokens, versus non object words. Here, the non object words are randomly selected to be the same number of token indices as the object words. We again see that while there is some minor information bleed, the bulk of spatial ID information lies in object word tokens.

C.3 STEERING EFFECTS ON ORTHOGONAL DIRECTIONS (X VS. Y), (TIME VS. X)

In Fig. 4, we show the results of horizontal steering on “left” vs. “right” beliefs, and vertical steering on “above” and “below”. To verify that steering directions can be decoupled, we perform the same steering and observe affects on beliefs of orthogonal directions. We show results of this preliminary analysis on LLaVA. Fig. A16 shows these orthogonal effects. Spatial IDs that are equivalent in the y coordinate but changing in x coordinate do not change beliefs in “above” or “below”. Similarly, static x coordinates with a changing y coordinate in spatial IDs has no effect on model belief about “left” and “right”.

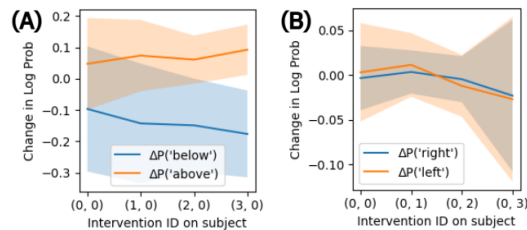


Figure A16: Steering effects of horizontal vectors on vertical beliefs (A) and vertical vectors on horizontal beliefs (B) in LLaVA.

Further, we check orthogonality between the space dimension and temporal dimension in video models. Fig. A17 shows a spatiotemporal ID grid from L11-14 on LLaVA-Video. The IDs are from videos where the object was in one of 8 frames (temporal change), and in one of 3 locations (spatial change). The experimental setup was minimal due to compute limitations. But even in this minimal setting, we see that the spatial and temporal axes are well separated.

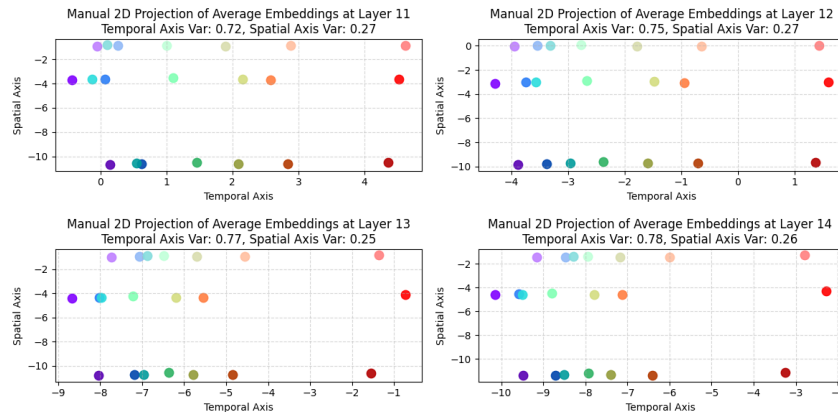


Figure A17: Spatiotemporal ID grid, where y axis is space and x axis is time.

D ABLATIONS

D.1 SCALING ANALYSIS FOR SPATIOTEMPORAL ID EXTRACTION

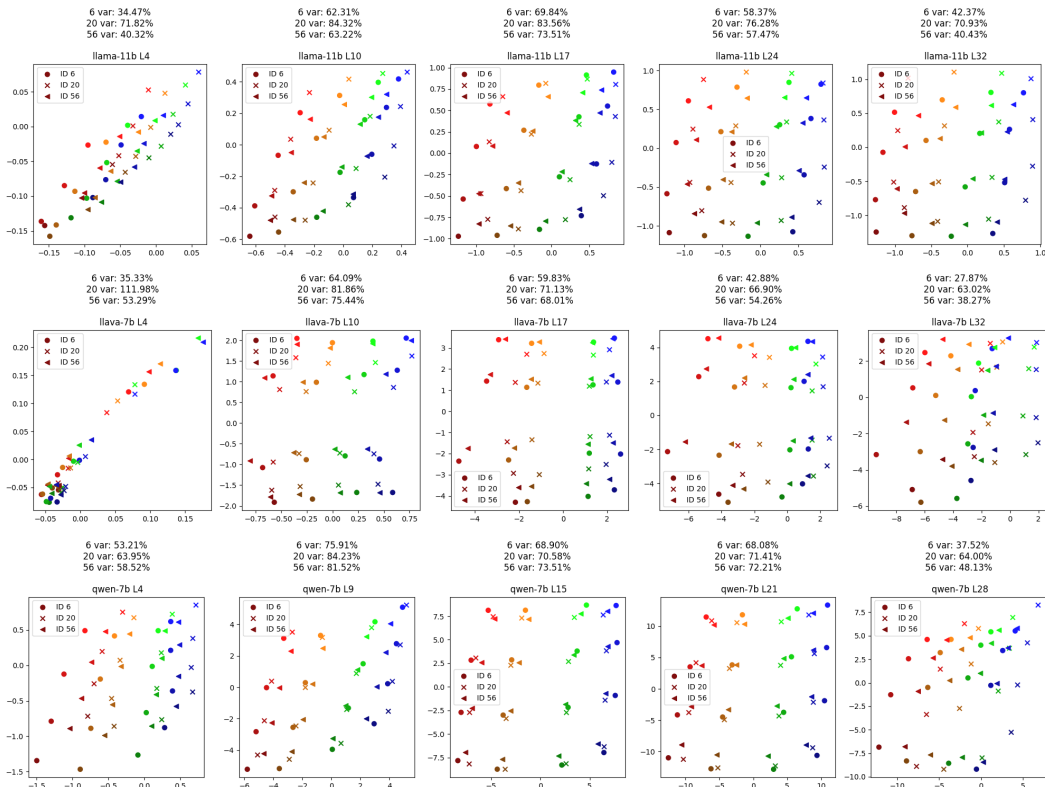


Figure A18: Extracting IDs with 6, 20, and 56 object pair images.

The projection axes are from the 56 object pair case. At intermediate layers, where we expect spatial IDs to be most crucial, we see a tight color-wise clustering, indicating spatial IDs extracted from various numbers of objects still converge. The variance explained by the spatial axes for all spatial ID extraction cases is $\gtrsim 50\%$, showing even at as little as 9 object pairs, we can extract good spatial IDs.

D.2 VARYING PROMPT WORDING AND OBJECT SIZES DURING EXTRACTION

Varying prompt wording. In this work, we use a spatial query in the form “Is the x to the left or right of the y?” to extract spatial IDs from object words. To verify that the choice of prompt does not matter, and that information about spatial location of objects flows into the word activation regardless, we extract spatial IDs from a *plain* prompt in the form “Is there an x or y in the image?”. In Fig. A19 we show the results of this extraction. We see that regardless of the input query formatting, spatial information can be extracted from the object words at intermediate.

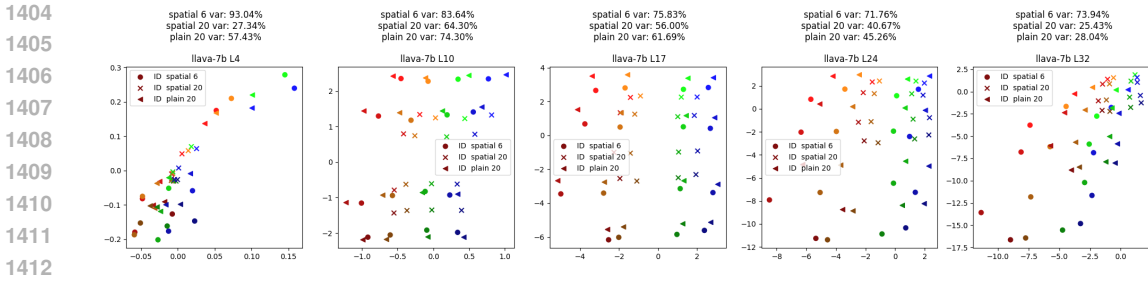


Figure A19: Plain prompts and spatial prompts projected onto spatial axes created from spatial prompts. Colors exhibit tight clustering.

Varying Image Sizes. To test that spatial IDs are roughly agnostic to object size, we extract spatial IDs from images where the object is 80px in diameter, 128px, and 176px, then project all extracted spatial IDs onto spatial axes created only from the medium sized object case. The result is shown in Fig A20. While the variance explained by the spatial axes drops by 10~20%, the spital IDs extracted from different sized objects still exhibit strong in-color clustering, and $\gtrsim 50\%$ of variance are explained by the spatial axes.

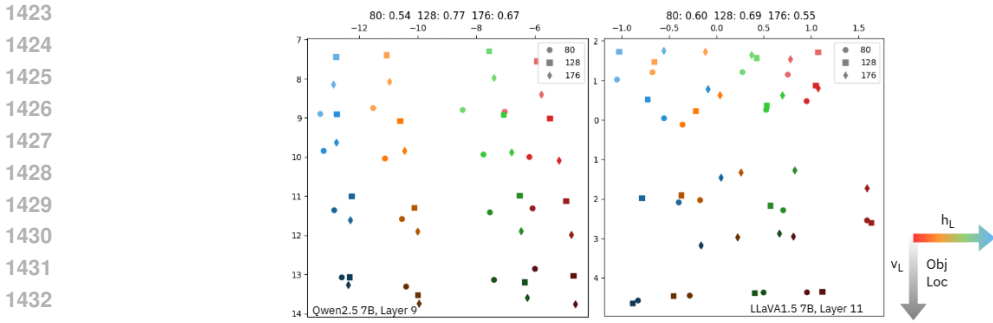


Figure A20: Spatial ID grids for Qwen and LLaVA, extracted from multiple object sizes. Circles are IDs extracted from images where object size was 80px in diameter, square is 128px, and diamond is 176px. On the top row, is the variance explained by the spatial axes for each size case.

E THEORETICAL ANALYSIS OF SPATIAL IDS

E.1 INFORMAL PROOF FOR SPATIAL ID EMERGENCE

Proposition: *Universal spatial IDs arises in any VLM using positional encoding, per self attention (Vaswani et al., 2017).*

Preliminaries. Consider a VLM layer with one attention head. Let the input sequence contain projected visual tokens $\{x_p\}_{p \in \mathcal{P}}$, where each patch index $p = (i, j)$ lies on an $m \times m$ grid, and text tokens including an object token o (as in prompts “Is there an o ?”). Define queries, keys, and values $q_o = W_Q r_o$, $k_p = W_K x_p$, $v_p = W_V x_p$, and the standard residual update $r_o \leftarrow r_o + W_{out} \sum_p \alpha_{o \leftarrow p} v_p$ with $\alpha_{o \leftarrow p} = \text{softmax}_p(q_o^\top k_p / \sqrt{d})$.

We make two very weak assumptions.

(1) First, we approximate that each patch vector decomposes as

$$x_p = s_p + P \psi(p) + \varepsilon_p,$$

where s_p encodes content (semantics), $\psi(p) \in \mathbb{R}^{d_\psi}$ is a shared positional basis (e.g., learned 2D embeddings or RoPE-induced features), P maps positional features into model space, and ε_p is some small deviation. In practice, explicit positional encoding is appended in autoregressive VLMs, so this assumption is explicitly true. In §E.2 we show empirically that positional encodings of VLMs linearly explain spatial IDs.

(2) We also assume that at a patch level, objectness is still encoded such that for images where a visual instance of the object word o occurs at a unique patch $p^* = (i, j)$, the attention kernel is peaked at p^* . In other words, $q_o^\top k_{p^*} \gg q_o^\top k_p$ for $p \neq p^*$, so that $\alpha_{o \leftarrow p^*} \approx 1$. Again, this is almost always true in practice, as modality alignment is encouraged during training.

Proof. Write the value at patch p using the decomposition:

$$v_p = W_V x_p = W_V s_p + W_V P \psi(p) + W_V \epsilon_p \quad (12)$$

Under assumption (2), the attention update to the object token is

$$\delta r_o = W_{out} \sum_p \alpha_{o \leftarrow p} v_p \approx W_{out} W_V x_{p^*} \quad (13)$$

Then we can rewrite Eq.3 as:

$$\begin{aligned} \Delta_L^{(o)}(p^*) &= r_{o,p^*} - \overline{r_{o,p}} \\ &= \left(r_o + W_{out} W_V (s_{p^*} + P \psi(p^*) + \epsilon_{p^*}) \right) - \left(r_o + W_{out} W_V (s_p + P \psi(p) + \epsilon_p) \right) \\ &= W_{out} W_V P \left(s_{p^*} - s_p + \psi(p^*) - \psi(p) + \epsilon_{p^*} - \epsilon_p \right) \end{aligned} \quad (14)$$

Note that $s_{(o,p^*)} = s_{(o,p)}$ for any p , for the first initial text embedding. Therefore, we can reduce Eq. 14 into:

$$\Delta_L^{(o)}(p^*) = \Delta_L^{(o)}(i, j) \simeq W_{out} W_V P \left(\psi(i, j) - \overline{\psi(p)} \right) \quad (15)$$

This expression is independent of o except through the common matrix $W_{out} U$, so averaging over objects leaves it unchanged. (In practice, we perform the averaging to reduce background noise.)

Notice that $W_{out} W_V P = M$ is fixed for some frozen network, and independent of location. Hence the centered attention update to the object token recovers a fixed linear ID of a shared positional basis, i.e., a universal spatial ID. The implications of the emergence of these intermediate IDs is that a shared spatial vocabulary need only be aligned with their respective positional basis vectors to perform “reasoning”. Let z_o be the residual stream at the object token after the update, and let W_{vocab} be the (approximately linear) readout to logits. Then

$$\ell(\text{LEFT}) - \ell(\text{RIGHT}) \approx (w_{\text{LEFT}} - w_{\text{RIGHT}})^\top \Delta_L(i, j) \approx (w_{\text{LEFT}} - w_{\text{RIGHT}})^\top M (\psi(i, j) - \mu_\psi) \quad (16)$$

so if $(w_{\text{LEFT}} - w_{\text{RIGHT}})^\top M$ aligns with the x -coordinate component of ψ , the model correctly predicts spatial words.

Multi-head and multi-layer accumulation. For H heads, $M = \sum_{h=1}^H W_{out}^{(h)} W_V^{(h)} P^{(h)}$; across layers, the contribution composes linearly in the residual stream. The “alignment band” in our experiments corresponds to layers where $\|M\|$ (or its projection onto the readout) is largest.

E.2 EMPIRICAL RELATIONSHIP BETWEEN POSITIONAL ENCODING AND SPATIAL IDS

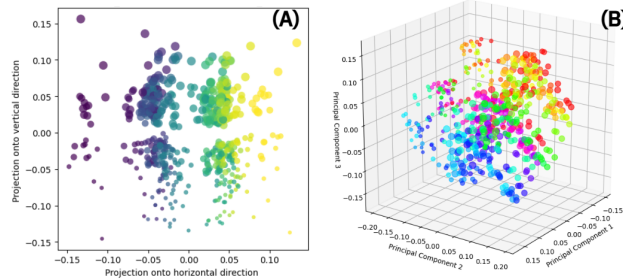


Figure A21: LLaVAPositional Encodings

Fig. A21 shows the patch level positional encodings from LLaVA(which uses the CLIP ViT-L/14 image encoder) projected onto 2 computed spatial axes or 3 principal components. The learned positional encoding vectors clearly have a linear structure, and with reduction in dimension are a linear transformation of the spatial ID grids we extract in §2.2. For a model like Qwen, which starts with fixed Rotary Positional Encodings (RoPE) that are not learned, this separable structure is innate. Previous work has shown that positional encoding in vision encoders continues to be linearly recoverable at penultimate activations (Ren et al., 2023). We are interested in whether this structure is linearly recoverable in a downstream LLM, in the form of spatial IDs, to support §E.1. We show that for the models studied, there exist low rank linear mappings from positional encodings to spatial IDs.

Setup. Let $X \in \mathbb{R}^{N \times d}$ be positional encodings for some model and $Y \in \mathbb{R}^{N \times M}$ be the spatial IDs extracted. To find their linear relationship, we simply must solve for $W \in \mathbb{R}^{d \times M}$ in $Y \approx XW$.

The least-squares solution is obtained with the Moore–Penrose pseudoinverse as $W^* = X^+Y$. To impose a rank constraint r , we compute the truncated singular value decomposition $X = U\Sigma V^\top$ and keep only the top r singular values Σ_r . Then the rank- r solution is

$$W_r = V_r \Sigma_r^{-1} U_r^\top Y \quad (17)$$

The in-sample fit can be quantified by the coefficient of determination:

$$R_r^2 = 1 - \frac{\|Y - XW_r\|_F^2}{\|Y - \bar{Y}\|_F^2}, \quad (18)$$

where \bar{Y} is the column-wise mean of Y .

For models like LLaVA and LLaMA, we acquire X by taking the learned positional embeddings. For models that use RoPE (which encodes position through complex rotations applied to query-key pairs) such as Qwen, we need an additional step to extract X . Specifically, we can form a RoPE design matrix from the sinusoidal basis functions underlying these rotations and perform the same reduced-rank regression to the extracted spatial IDs Y . Each position $p \in \{0, \dots, N-1\}$ is mapped to sinusoidal features at different frequencies. Let the hidden dimension be d , with frequencies

$$\theta_i = 10000^{-\frac{2i}{d}}, \quad i = 0, \dots, \frac{d}{2} - 1.$$

The RoPE design matrix $X_{RoPE} \in \mathbb{R}^{N \times d}$ is then

$$X_{RoPE}(p) = [\cos(\theta_0 p), \sin(\theta_0 p), \cos(\theta_1 p), \sin(\theta_1 p), \dots, \cos(\theta_{d/2-1} p), \sin(\theta_{d/2-1} p)], \quad (19)$$

with each row of Φ corresponding to a position p .

Results. We find that a weight matrix of rank 3 linearly relates the positional encoding matrix to the spatial IDs of a model with $R^2 \geq 0.85$. The three independent weight vectors likely correspond to horizontal, vertical, and radial axes, meaning such structure is preserved in the spatial IDs with high fidelity. Results are shown in Table 2.

Model	Rank-2 R^2	Rank-3 R^2
LLaVA1.5-7B	0.458	0.854
LLaMA3.2VL-11B	0.610	0.869
Qwen2.5VL-7B	0.605	0.903

Table 2: R^2 from low rank W

F LLM USAGE DISCLOSURE

GPT-4 and GPT-5 were used in the process of occasionally coding experiments and editing paper wording.