

# Linear Mechanisms for Spatiotemporal Reasoning in Vision Language Models

Anonymous CVPR submission

Paper ID \*\*\*\*

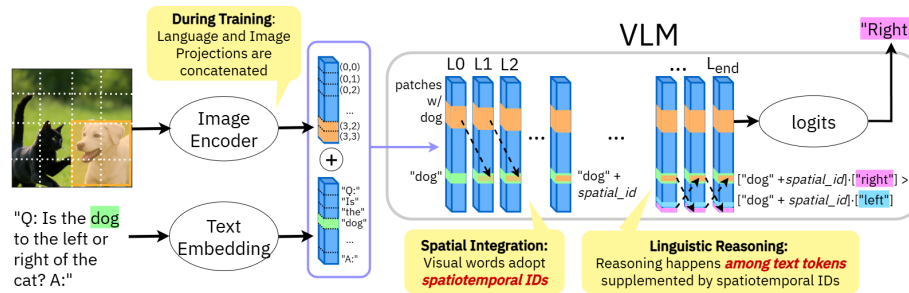


Figure 1. **Hypothesis for spatiotemporal visual reasoning.** The VLM linearly binds spatiotemporal localization to object word activations in early layers. Subsequent linguistic reasoning about the object is informed by its location in space and time per the spatiotemporal ID.

## Abstract

001 *Spatio-temporal reasoning is a remarkable capability of*  
 002 *Vision Language Models (VLMs), but the underlying mech-*  
 003 *anisms of such abilities remain largely opaque. We postu-*  
 004 *late that visual/geometrical and textual representations of*  
 005 *spatial structure must be combined at some point in VLM*  
 006 *computations. We search for such confluence, and ask*  
 007 *whether the identified representation can causally explain*  
 008 *aspects of input-output model behavior through a linear*  
 009 *model. We show empirically that VLMs encode object loca-*  
 010 *tions by linearly binding spatial IDs to textual activations,*  
 011 *then perform reasoning via language tokens. Through rig-*  
 012 *orous causal interventions we demonstrate that these IDs,*  
 013 *which are ubiquitous across the model, can systematically*  
 014 *mediate model beliefs at intermediate VLM layers. Addi-*  
 015 *tionally, we find that spatial IDs serve as a diagnostic tool*  
 016 *for identifying limitations in existing VLMs, and as a valu-*  
 017 *able learning signal. We extend our analysis to video VLMs*  
 018 *and identify an analogous linear temporal ID mechanism.*  
 019 *By characterizing our proposed spatiotemporal ID mecha-*  
 020 *nism, we elucidate a previously underexplored internal re-*  
 021 *asoning process in VLMs, toward improved interpretability*  
 022 *and the principled design of more aligned and capable mod-*  
 023 *els.*

## 1. Introduction

Reasoning about visual input with textual queries is a key  
 challenge behind vision-language models (VLMs). Con-  
 sider 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 cat-  
 egorical decision. Though complex capabilities in spatial or  
 temporal reasoning are still far from being fully understood  
 or reliably engineered [5, 32, 34], 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 vari-  
 ous structural properties emerge in VLM internals during  
 VQA [17, 27, 42]. Relatedly, mechanistic interpretability in  
 LLMs has uncovered linear circuits for relational linguistic  
 reasoning [12, 15, 28]. 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?*

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 fami-  
 lies using this paradigm are LLaVA[24], LLaMA[9], Qwen  
 [3], InternVL [6], and Gemma [33]. A growing body of  
 work aims to improve spatial reasoning capacities in VLMs  
 [4, 11] and temporal reasoning in video models [22, 37].

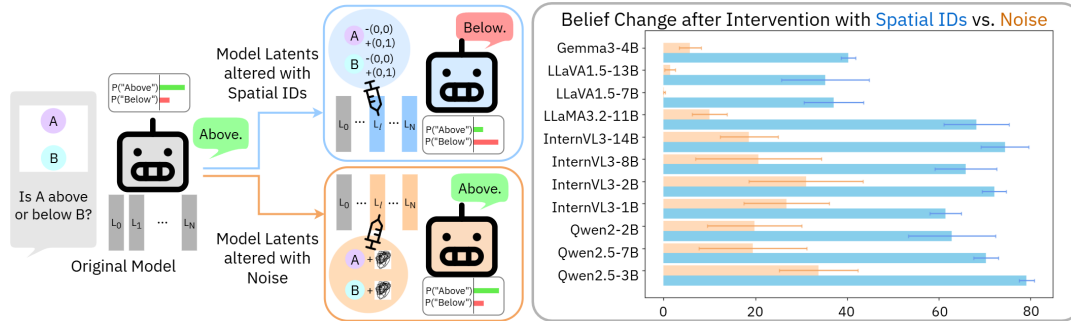


Figure 2. **Results from Targeted Intervention** (§3). Median binary belief swap due to spatial ID steering is 64.4%, and 29.5% for noise. Spatial IDs have 43.6% above-chance influence on average. We conclude that spatial IDs mediate models’ beliefs about objects’ locations in space.

053 Identification of the internal mechanism by which SoTA  
 054 VLMs do spatial VQA can empower engineers to identify  
 055 current architectural components leading to VQA failure  
 056 modes in 3D reasoning or simple VQA. To this end, we ask:  
 057 **Q2.** *Given our linear model of spatial reasoning in model*  
 058 *activations, how do we use it to understand and improve*  
 059 *SoTA VLMs?*

060 Similar training paradigms to image-based VLMs  
 061 yield video models such as LLaVA-Video[43],  
 062 VideoLLaMA3[40], and Qwen2.5 [3], among others.  
 063 Given our theory for the mechanisms underlying spatial  
 064 reasoning in VLMs, we ask: **Q3.** *Do video models utilize*  
 065 *analogous linear mechanisms for temporal reasoning?*

066 To address these questions, we conduct a mechanis-  
 067 tic analysis of autoregressive VLMs and construct a linear  
 068 model for spatiotemporal reasoning. We show that VLMs  
 069 decompose a visual reasoning task by first binding spatial  
 070 information about visual objects to object word activations,  
 071 in the form of linear components we term *spatial IDs*,  
 072 answering Q1 (Fig. 1). We then extract these IDs and demon-  
 073 strate their mediative capacity on model output through tar-  
 074 getted belief steering in text activations (Fig. 2). We fur-  
 075 ther find that spatial IDs provide insight on VLMs’ struggle  
 076 with depth reasoning, and incorrect spatial IDs as a result of  
 077 weak vision encoder or poor modality integration leads to  
 078 failures in LLaVA and LLaMA. This answers Q2. Finally,  
 079 we show that temporal IDs similarly mediate video models,  
 080 answering Q3. In summary, our novel contributions are:

- 081 • **Spatial ID Model Formulation:** We propose a linear  
 082 model of spatial reasoning in VLMs, called *spatial IDs*.  
 083 These are text-anchored latent structures that bind visual  
 084 elements to object tokens thus enabling linguistic reason-  
 085 ing about space (§2.1). We empirically extract them from  
 086 SoTA VLMs for characterization (§2.2).
- 087 • **Analytical and Empirical Proof of Causality:** We show  
 088 model belief can be manipulated by perturbing only the  
 089 spatial IDs, demonstrating their causal role in reasoning  
 090 (§3), and provide theoretical intuition for the emergence  
 091 of spatial IDs in VLMs (§2.3).

- **SoTA VLM Analysis with Spatial IDs:** Through tar-  
 getted intervention, we identify limitations in depth  
 expression (§4.1) and systematic failure modes in  
 LLaMA/LLaVA (§4.2).
- **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.

### 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 counter-

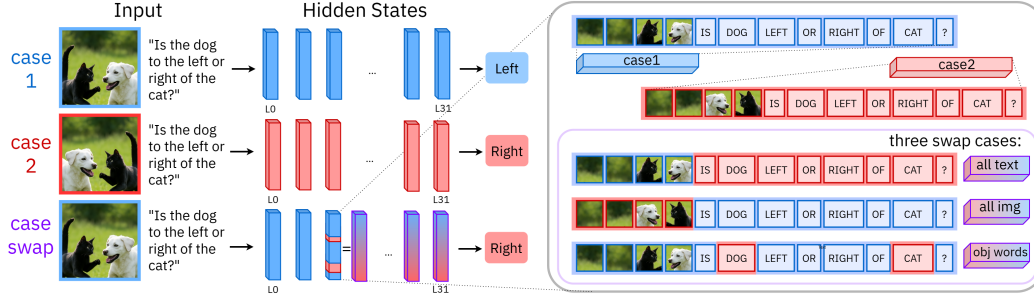


Figure 3. **Mirror swapping experiment** (§2.1). Activations from case 1 and 2 are partially swapped at a select layer, in one of three arrangements. Computations continue normally after this point.

128 part  $y_L$ . The modified representation  $\tilde{x}_L$  is passed through  
 129 the remaining layers. We conduct interventions with three  
 130 variants of  $Q$ : (1) all text tokens (2) all image patches (3)  
 131 object-word tokens only. If information critical to spatial  
 132 reasoning is concentrated in any of these, the model’s belief  
 133 will change when that region is overwritten. As a control,  
 134 we concurrently perform “attribute swapping”, which fol-  
 135 lows the same steps but instead of mirroring the input image  
 136 for the intervention case, changes its colors. The interven-  
 137 tion procedure is visualized in Fig. 3 and formally defined  
 138 in Alg. 1. Further implementation specifics are deferred to  
 Appendix §A.1.

**Algorithm 1** Swapping Intermediate Activations in Mirrored Images

$$\begin{aligned}
 x_L &\leftarrow f_L \circ \dots \circ f_1(x) \\
 y_L &\leftarrow f_L \circ \dots \circ f_1(y) &> x, y: [\text{seq\_dim}, \text{embed\_dim}] \\
 \tilde{x}_L &\leftarrow x_L[\tilde{Q}] + y_L[Q] &> \tilde{x}_L: [\text{seq\_dim}, \text{emb\_dim}] \\
 \tilde{x}_{\text{out}, L} &\leftarrow f_{L_{\text{max}}} \circ \dots \circ f_{L+1}(\tilde{x}_L) \\
 y_{\text{out}} &\leftarrow f_{L_{\text{max}}} \circ \dots \circ f_{L+1}(y_L) &> P_{\tilde{x}_{\text{out}, L}}: [1]
 \end{aligned}$$

139 Here,  $Q$  denotes the selected indices in the input se-  
 140 quence to swap, and  $\tilde{Q}$  is all other indices. We use the  
 141 COCO-SPATIAL benchmark [18] for the mirrored images,  
 142 which is a curated subset of COCO [23] annotated with spa-  
 143 tial language. To quantify belief shift caused by the in-  
 144 tervention, we compute the fraction of the mirror-induced  
 145 change that can be attributed to the swapped activations at  
 146 layer  $L$ . For the ground truth logit “GT”, this quantity is  
 147 derived as:  
 148

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

150 **Results from Mirror Swapping** are shown in Fig. 4A.  
 151 Through mirror swapping, we observe a *layer-specific ef-*  
 152 *fect* for intervention effect across modalities. Intervening  
 153 on visual patch tokens has a strong effect in early layers  
 154 but fades with depth. Conversely, interventions on text to-  
 155 kens increasingly affect final outputs in later layers. This  
 156 is corroborated by observations that middle layers have a  
 157 modality switching effect in VLMs [17]. 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 mir-  
 ror 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 equiva-  
 lent 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  $\sim 20 \times$  that of the original be-  
 lief difference.

These results suggest that VLMs extract and encode spa-  
 tial facts from the image into object word tokens’ activa-  
 tions, 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 ac-  
 tivations, 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 in-  
 deed linear.

**Extraction Preliminaries.** We first set up some for-  
 malisms 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 posi-  
 tioned 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 de-  
 fine  $\phi_L(o; I_{(i,j)}^{(o)}, T^{(o)}) \in \mathbb{R}^d$  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:

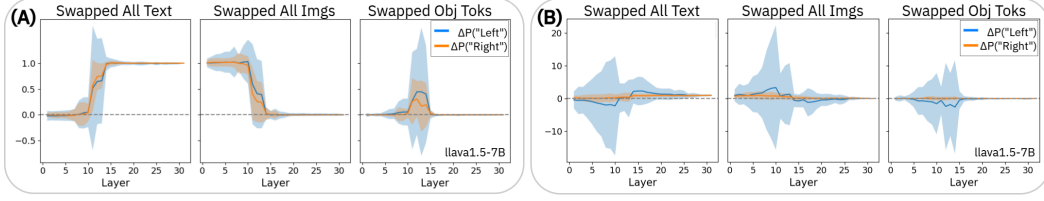


Figure 4. **Ratio change in log probability for logits “left” and “right” from mirror swap (A) and attribute swap (B) interventions.** (A) shows distinct binary belief swaps, where text tokens have an influence after middle layers. Image patches stop having an influence after that point, and object word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

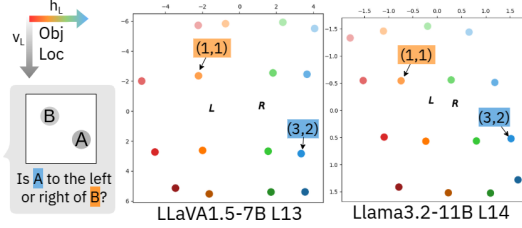


Figure 5. **Spatial IDs in a grid.** Color and saturation of markers represent the location of the object when spatial ID was extracted. x and y axes are coefficients of ID projections onto  $h_L$  and  $v_L$ . L, R represent “left”, “right” textual activations.

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

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 OBJVERSE [8] 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 and counterfactual studies confirming that spatial IDs reside in object

words, and spatial axes are orthogonal. 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:

$$r_o \leftarrow r_o + W_{\text{out}} \sum_{p \in P} \alpha_{o \leftarrow p} v_p, \quad (6)$$

$$\alpha_{o \leftarrow p} \propto \exp\left(\frac{q_o^\top k_p}{\sqrt{d}}\right), \quad v_p = W_V x_p.$$

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  $\varepsilon_p$  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:

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

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  $\psi$ , the model prefers “left.” Empirically, a low-rank linear fit from positional encodings  $\psi$  to spatial IDs  $\Delta_L$  explains

258 most variance (e.g. rank-3 gives  $R^2 \gtrsim 0.85$ , see §E.2, Ta-  
259 ble 1). A more detailed derivation for  $\Delta_L(i, j)$  for the mul-  
260 ti-head case is shown in Appendix §E.1. This is a particu-  
261 larly simplified setting, and real reasoning circuits in VLMs  
262 will involve a lot more noise and nonlinearities. The main  
263 takeaway is that VLM designs like Fig. 1 encourage models  
264 to endow text tokens with visual information, followed by  
265 linguistic reasoning. This information transfer, in its most  
266 simplified linear form, is in the form of spatial IDs.

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

### 276 3. Spatial IDs Mediate Model Beliefs

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

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

**Algorithm 2** Intervention at Layer  $L$  via Residual Modifi-  
cation

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

289 Here we scale the norm of  $\Delta_L(i, j)$  to be  $\alpha|x_L[q]|$ , and  
290  $\tilde{\Delta}_L(i, j) = \Delta_L(m-i-1, j)$ . This approximately preserves  
291 the norm of  $x_L$ .  $\alpha = 5$  is some scaling constant set after  
292 grid searching for stable intervention. We take 100 COCO  
293 images where one object is to the left or right of another, per  
294 labels from COCO-SPATIAL, and ask queries of the form “Is  
295  $x$  to the left/right of  $y$ ?”. We measure the log probability of  
296 “left” and “right” tokens in the final output logits to assess  
297 steering effects.

298 **Results from Arbitrary Steering.** Fig. 6 shows the ef-  
299 fects of model belief steering on real images and videos.  
300 Fig. 6A shows that steering impact is greatest at modality  
301 alignment layers as expected per the mirror swapping anal-

302 ysis, and Fig 6B shows that intervening with the rightmost  
303 spatial ID largely enhances model belief that the object is  
304 to the right, and vice versa for the leftmost ID for leftward  
305 belief. The y axes show changes in log probability for those  
306 binary directions for the whole dataset, and x axes show the  
307 different ID locations. Regardless of whether the answer to  
308 the original query was “left” or “right”, subplot trends are  
309 the same.

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

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

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

### 345 4. Spatial IDs for Understanding and Improv- 346 ing Image VLMs

347 With the existence and causal nature of spatial IDs es-  
348 tablished, we explore two ways to leverage them towards  
349 stronger VLMs. First, we aim to understand why 3D rea-  
350 soning fails in SoTA VLMs. Second, we use spatial IDs to  
351 diagnose architectural bottlenecks of SoTA VLMs in VQA.

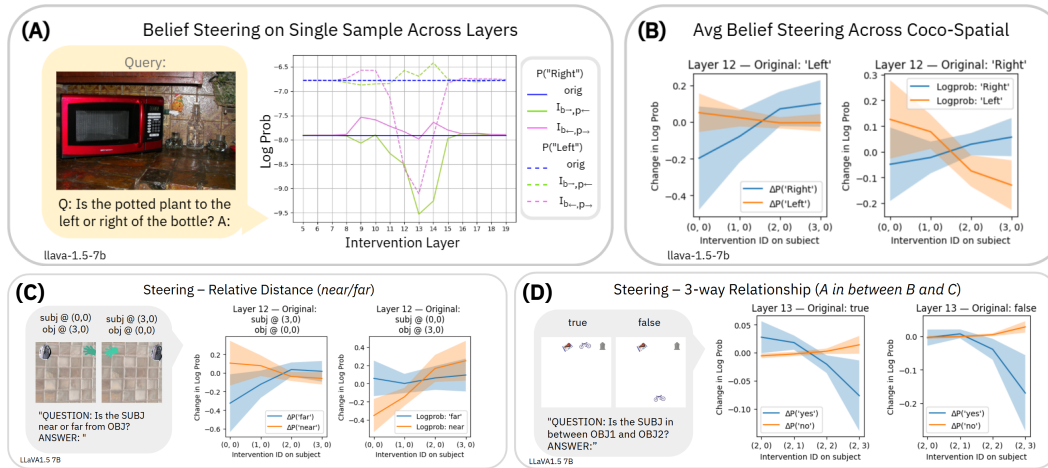


Figure 6. **Effect of spatial steering on real images** on one sample across different intervention layers (A) and across a dataset for one layer (B). In (A), dotted and solid lines indicate answer probabilities for “left” or “right”. Different colors indicate no intervention (blue), steering the bottle 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”.

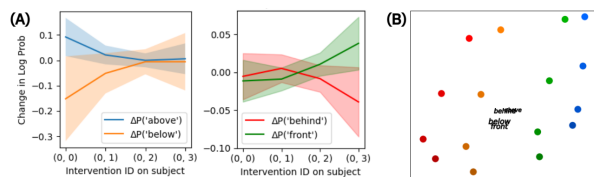


Figure 7. **Depth and height are strongly correlated in LLaVA.** (A) Steering results for IDs varying in y-dim and their impact on beliefs about height or depth. (B) Projection of spatial words onto a spatial ID grid. Embeddings for “above”/“front” and “behind”/“below” are nearly identical.

## 352 4.1. Depth Representation in Image VLMs

353 Spatial IDs suggest that VLMs represent visual space within  
354 a 2D grid. What might this mean for depth? We hypothe-  
355 size that the language model must reason about depth re-  
356 lated queries using the 2D localization in context. To ver-  
357 ify whether this is the case, we look at the resulting belief  
358 changes in the depth axis when the LLaVA 1.5 7B model is  
359 steered with spatial IDs varying in height. Fig. 7 shows the  
360 results. The same spatial IDs increasing the likelihood for  
361 “above” and decreasing “below”, also drive up “front” and  
362 drive down “behind” in LLaVA.

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

## 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 fail-  
ure 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 sam-  
ples. In this section, we perform two experiments to identify  
the architectural component which causes for the distribu-  
tion 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 incor-  
rect 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 sensi-  
tivity 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 cor-  
related to wrong answers, but still have some model inaccura-  
cies be resultant from factors other than spatial IDs, such as  
erroneous priors during LM readout [20, 30]. 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 truth (g.t.), we compute the g.t. spatial ID by projecting the word activation onto the spatial axes:

$$\Delta_L^{(o)}(i, j)_{ext} \approx VV^T \phi_L(o; I_{(i,j)}^{(o)}, T^{(o)}), \quad (9)$$

$$V = [v_L, h_L]$$

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

$$\text{ID deviation margin} = \epsilon_{ext} - \epsilon_{gt},$$

$$\text{where } \epsilon_{gt} = i_{gt}^{(o)} - i_{gt}^{(\tilde{o})}, \quad (10)$$

$$\epsilon_{ext} = i_{ext}^{(o)} - i_{ext}^{(\tilde{o})}$$

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  $\epsilon_{ext} - \epsilon_{gt}$ . The red histogram represents the density of ID deviations for  $\mathbf{K}_{wrong}$ , and the blue histogram shows the same for  $\mathbf{K}_{correct}$ . The red distribution is visibly skewed to the negatives compared to the blue. Quantitatively, we perform the Mann-Whitney U test [25] 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  $\mathbf{K}_{correct}$  are more sensitive to masking the image raw input at the g.t. location of  $o$ , while beliefs on  $\mathbf{K}_{wrong}$  change more with masking elsewhere, we can conclude that

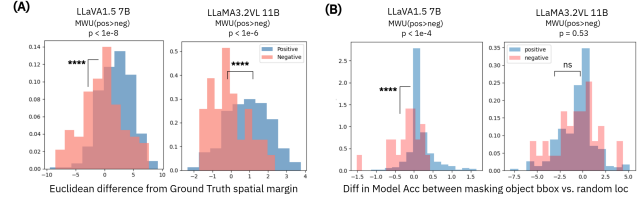


Figure 8. **Contrasting density histograms show incorrect spatial IDs drive bad predictions.** (A) shows deviation of model spatial IDs from g.t., and (B) the difference in model accuracy when masking objects vs. random locations in images. Histograms are samples VLMs got right (blue) or wrong (red).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

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 needed capacity.

We design an obfuscation paradigm inspired by methods like D-RISE [29], 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  $\tilde{o}$ . We then measure model belief change when masking the object vs. elsewhere:

$$\text{bbox sensitivity} =$$

$$(P(\text{“GT”}) - P(\text{“GT”}|\text{mask } o)) \quad (11)$$

$$- \left( P(\text{“GT”}) - \min_r [P(\text{“GT”}|\text{mask } r), r \in R] \right)$$

**Image Masking Results.** Fig. 8B shows overlaid histograms for bounding box masking sensitivities of  $\mathbf{K}_{correct}$  and  $\mathbf{K}_{wrong}$ . 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  $\mathbf{K}_{wrong}$  is shifted more negative than  $\mathbf{K}_{correct}$ , indicating its vision encoder fails at object detection when it answers incorrectly. In contrast,  $\mathbf{K}_{wrong}, \mathbf{K}_{correct}$  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 fine-grained 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 re-

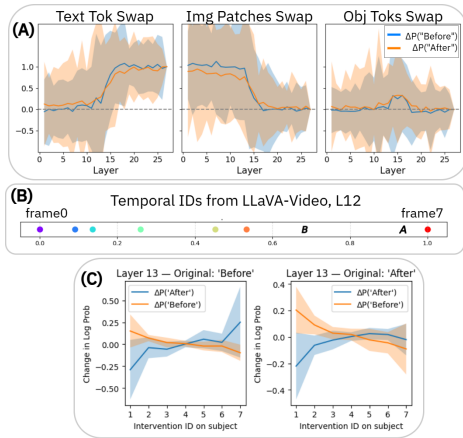


Figure 9. **Temporal ID Results.** Mirror Swapping on videos (A), Temporal ID grid (B), and temporal ID steering on model beliefs (C) with LLaVA-Video

490 sources are scarce, or provide intuition for future model de-  
491 signs.

## 492 5. Temporal IDs in Video Models

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

### 501 5.1. Mirroring, Extracting, and Steering across the 502 Temporal Axis

503 **Temporal Mirror Swapping.** We validate that there exist  
504 modality alignment layers with object-level visual informa-  
505 tion transfer in video models. For mirrored videos, we take  
506 the Scene\_QA subset of MVBENCH [21] and swap the order  
507 of frames from back to front. Following Alg. 1, we  
508 show results for swapping text tokens, image patches, and  
509 object words in Fig. 9A. While the error bound is noisier  
510 than spatial LLaVA, likely as LLaVA-Video follows re-  
511 sponse formats less well, we see the expected bump around  
512 middle layers for crossmodal integration.

513 **Temporal ID Extraction.** Derivation of temporal IDs  
514 and the temporal vector  $t_L$  follows Eq. 2 - 5, with synthetic  
515 8-frame videos of OBJVERSE renders. Results are shown  
516 in Fig. 9B. We again see that the text activation for “be-  
517 fore” projects closer to earlier frames, than the activation  
518 for “after”.

519 **Causality of Temporal IDs.** Finally, to confirm control-  
520 lability with arbitrary temporal IDs, we perform the steering  
521 experiment per Alg. 2 on MVBENCH videos. Results are  
522 shown in Fig. 9C. On these real, naturalistic videos, we see  
523 that later temporal IDs steer the model belief towards “af-  
524 ter”, and earlier IDs towards “before”, as expected.

## 525 5.2. Emergence of Temporal IDs

526 Fig. 9 shows summary results on LLaVA-Video, but we in-  
527 clude temporal IDs from VideoLLaMA3 and Qwen2.5 in  
528 Appendix §B.2. LLaVA-Video and VideoLLaMA3 use text-  
529 ual description of the video length and number of frames  
530 to indicate timestamps preceding the visual input, while  
531 Qwen uses explicit MRoPE time IDs. This suggests that  
532 spatiotemporal IDs can emerge without explicit positional  
533 encoding, beyond the simple mechanism derived in Eq. 7.

## 534 6. Related Work

535 Mechanistic interpretability is a growing field uncovering  
536 the inner workings of large models, popularizing techniques  
537 such as circuit tracing [2, 10], Sparse Autoencoders [7], lin-  
538 ear probing [1], and others. The Linear Representation Hy-  
539 pothesis posits that concepts are linearly encoded in LLM  
540 latents [28], and activation patching supports that linear  
541 changes in activations drive model belief [26, 41]. Internal  
542 in-context reasoning mechanisms such as linear *binding IDs*  
543 [12, 13] have been identified in LLMs, along with other ev-  
544 idence for linear multi-hop reasoning [38], in-context task  
545 vectors [14] and linear relational embeddings [15] during  
546 reasoning.

547 Linearity of embeddings have also been discovered in  
548 VLM latent spaces [16, 35] to some degree. Previous work  
549 showed that VLMs separate VQA into image-focused then  
550 text-focused stages [17], and others have extended LLM  
551 interpretability techniques like logit lens [27] or attention  
552 tracking [39, 42] to VLMs to unearth internal circuits. In  
553 our work, we mechanistically capture spatiotemporal infor-  
554 mation flow from image patches to text tokens in VLMs, via  
555 the spatial ID mechanism.

## 556 7. Conclusion, Limitations, & Future Work

557 We propose spatiotemporal IDs as a linear model for visual  
558 reasoning about space and time in VLMs. With a series of  
559 causal analyses, we show these IDs can be obtained in many  
560 SoTA models, and that they closely mediate models’ beliefs  
561 about visual objects’ location in space and time. We further  
562 offer ways to extend this mechanistic insight to improv-  
563 ing existing VLMs. For tractability, our work is currently  
564 limited to analyses in simple spatial queries or appearance-  
565 based temporal queries. Experimental design for more com-  
566 plex, open-ended queries will enhance our understanding of  
567 how VLMs utilize rudimentary concepts like spatial IDs in  
568 more nuanced settings. Further, we only extract and steer  
569 models of sizes up to 14B parameters. Investigation into  
570 whether the spatial ID circuit plays a similarly prominent  
571 role in larger models will reveal whether VLMs of varying  
572 capacities follow analogous methods for visual reasoning,  
573 or employ distinct measures. Lastly, while we show several  
574 potential ways to leverage spatial IDs for VLM diagnostics  
575 or finetuning, future work could include expanded use cases  
576 such as explicit temporal guidance at large scale.

577

**References**

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

- [1] Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes. *arXiv preprint arXiv:1610.01644*, 2016. 8
- [2] 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 Jermyn, 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. 8
- [3] 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, 2025. arXiv:2502.13923 [cs]. 1, 2
- [4] 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*, pages 14455–14465, 2024. 1
- [5] 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. 1
- [6] 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, 2024. arXiv:2312.14238 [cs]. 1
- [7] 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. 8
- [8] 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*, pages 13142–13153, 2023. 4
- [9] 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*, pages arXiv–2407, 2024. 1
- [10] 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>. 8
- [11] 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, 2025. arXiv:2505.15879 [cs]. 1
- [12] Jiahai Feng and Jacob Steinhardt. How do Language Models Bind Entities in Context?, 2024. arXiv:2310.17191 [cs]. 1, 8
- [13] Jiahai Feng, Stuart Russell, and Jacob Steinhardt. Monitoring Latent World States in Language Models with Propositional Probes, 2024. arXiv:2406.19501 [cs]. 8
- [14] Roei Hendel, Mor Geva, and Amir Globerson. In-Context Learning Creates Task Vectors, 2023. arXiv:2310.15916 [cs]. 8
- [15] Evan Hernandez, Arnab Sen Sharma, Tal Haklay, Kevin Meng, Martin Wattenberg, Jacob Andreas, Yonatan Belinkov, and David Bau. Linearity of Relation Decoding in Transformer Language Models, 2024. arXiv:2308.09124 [cs]. 1, 8
- [16] Nick Jiang, Anish Kachinthaya, Suzie Petryk, and Yossi Gandelsman. Interpreting and Editing Vision-Language Representations to Mitigate Hallucinations, 2025. arXiv:2410.02762 [cs]. 8
- [17] Zhangqi Jiang, Junkai Chen, Beier Zhu, Tingjin Luo, Yankun Shen, and Xu Yang. Devils in Middle Layers of Large Vision-Language Models: Interpreting, Detecting and Mitigating Object Hallucinations via Attention Lens, 2025. arXiv:2411.16724 [cs]. 1, 3, 8
- [18] Amita Kamath, Jack Hessel, and Kai-Wei Chang. What’s” up” with vision-language models? investigating their struggle with spatial reasoning. *arXiv preprint arXiv:2310.19785*, 2023. 3
- [19] Raphi Kang, Yue Song, Georgia Gkioxari, and Pietro Perona. Is clip ideal? no. can we fix it? yes! *arXiv preprint arXiv:2503.08723*, 2025. 16

- 682 [20] Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li,  
683 Shijian Lu, Chunyan Miao, and Lidong Bing. Mitigating  
684 object hallucinations in large vision-language  
685 models through visual contrastive decoding. In *Proceedings of the IEEE/CVF Conference on Computer  
686 Vision and Pattern Recognition*, pages 13872–13882,  
687 2024. 6 734
- 689 [21] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi  
690 Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen,  
691 Ping Lou, Limin Wang, and Yu Qiao. MVBench:  
692 A Comprehensive Multi-modal Video Understanding  
693 Benchmark. In *2024 IEEE/CVF Conference on Com-  
694 puter Vision and Pattern Recognition (CVPR)*, pages  
695 22195–22206, Seattle, WA, USA, 2024. IEEE. 8 735
- 696 [22] Lei Li, Yuanxin Liu, Linli Yao, Peiyuan Zhang,  
697 Chenxin An, Lean Wang, Xu Sun, Lingpeng Kong,  
698 and Qi Liu. Temporal reasoning transfer from text to  
699 video. *arXiv preprint arXiv:2410.06166*, 2024. 1 736
- 700 [23] Tsung-Yi Lin, Michael Maire, Serge Belongie, James  
701 Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and  
702 C Lawrence Zitnick. Microsoft coco: Common ob-  
703 jects in context. In *European conference on computer  
704 vision*, pages 740–755. Springer, 2014. 3 737
- 705 [24] Haotian Liu, Chunyuan Li, Qingyang Wu, and  
706 Yong Jae Lee. Visual Instruction Tuning, 2023.  
707 *arXiv:2304.08485 [cs]*. 1 738
- 708 [25] Patrick E McKnight and Julius Najab. Mann-whitney  
709 u test. *The Corsini encyclopedia of psychology*, pages  
710 1–1, 2010. 7 739
- 711 [26] Kevin Meng, David Bau, Alex Andonian, and Yonatan  
712 Belinkov. Locating and editing factual associations in  
713 gpt. *Advances in neural information processing sys-  
714 tems*, 35:17359–17372, 2022. 8 740
- 715 [27] Clement Neo, Luke Ong, Philip Torr, Mor Geva,  
716 David Krueger, and Fazl Barez. Towards Interpreting  
717 Visual Information Processing in Vision-Language  
718 Models, 2024. *arXiv:2410.07149 [cs]*. 1, 8 741
- 719 [28] Kiho Park, Yo Joong Choe, and Victor Veitch. The  
720 Linear Representation Hypothesis and the Geometry  
721 of Large Language Models, 2024. *arXiv:2311.03658  
722 [cs]*. 1, 8 742
- 723 [29] Vitali Petsiuk, Rajiv Jain, Varun Manjunatha, Vlad I  
724 Morariu, Ashutosh Mehra, Vicente Ordonez, and Kate  
725 Saenko. Black-box explanation of object detectors via  
726 saliency maps. In *Proceedings of the IEEE/CVF con-  
727 ference on computer vision and pattern recognition*,  
728 pages 11443–11452, 2021. 7 743
- 729 [30] Sainandan Ramakrishnan, Aishwarya Agrawal, and  
730 Stefan Lee. Overcoming language priors in visual  
731 question answering with adversarial regularization.  
732 *Advances in neural information processing systems*,  
733 31, 2018. 6 744
- [31] Bin Ren, Yahui Liu, Yue Song, Wei Bi, Rita Cuc-  
chiara, Nicu Sebe, and Wei Wang. Masked jigsaw puzzle: A versatile position embedding for vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20382–20391, 2023. 29 745
- [32] Ilias Stogiannidis, Steven McDonagh, and Sotirios A Tsafaris. Mind the gap: Benchmarking spatial reasoning in vision-language models. *arXiv preprint arXiv:2503.19707*, 2025. 1 746
- [33] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024. 1 747
- [34] Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide shut? exploring the visual shortcomings of multimodal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9568–9578, 2024. 1 748
- [35] Matthew Trager, Pramuditha Perera, Luca Zancato, Alessandro Achille, Parminder Bhatia, and Stefano Soatto. Linear spaces of meanings: compositional structures in vision-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15395–15404, 2023. 8 749
- [36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. 27 750
- [37] Junbin Xiao, Angela Yao, Yicong Li, and Tat-Seng Chua. Can i trust your answer? visually grounded video question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13204–13214, 2024. 1 751
- [38] Sohee Yang, Elena Gribovskaya, Nora Kassner, Mor Geva, and Sebastian Riedel. Do Large Language Models Latently Perform Multi-Hop Reasoning?, 2024. *arXiv:2402.16837 [cs]*. 8 752
- [39] Zeping Yu and Sophia Ananiadou. Understanding Multimodal LLMs: the Mechanistic Interpretability of Llava in Visual Question Answering, 2025. *arXiv:2411.10950 [cs]*. 8 753
- [40] Boqiang Zhang, Kehan Li, Zesen Cheng, Zhiqiang Hu, Yuqian Yuan, Guanzheng Chen, Sicong Leng, Yuming Jiang, Hang Zhang, Xin Li, Peng Jin, Wenqi Zhang, Fan Wang, Lidong Bing, and Deli Zhao. VideoLLaMA 3: Frontier Multimodal Foundation Models for Image and Video Understanding, 2025. *arXiv:2501.13106 [cs]*. 2 754

- 787 [41] Fred Zhang and Neel Nanda. Towards best practices of  
788 activation patching in language models: Metrics and  
789 methods. *arXiv preprint arXiv:2309.16042*, 2023. 8
- 790 [42] Xiaofeng Zhang, Yihao Quan, Chen Shen, Xiaosong  
791 Yuan, Shaotian Yan, Liang Xie, Wenxiao Wang,  
792 Chaochen Gu, Hao Tang, and Jieping Ye. From Redundancy to Relevance: Information Flow in LVLMS  
793 Across Reasoning Tasks, 2024. *arXiv:2406.06579*  
794 [cs]. 1, 8
- 795
- 796 [43] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun  
797 Ma, Ziwei Liu, and Chunyuan Li. Video Instruction  
798 Tuning With Synthetic Data, 2024. *arXiv:2410.02713*  
799 [cs]. 2