
TEXTUAL STEERING VECTORS CAN IMPROVE VISUAL UNDERSTANDING IN MULTIMODAL LARGE LANGUAGE MODELS

006 **Anonymous authors**

007 Paper under double-blind review

ABSTRACT

Steering methods have emerged as effective tools for guiding large language models' behavior, yet multimodal large language models (MLLMs) lack comparable techniques due to recency and architectural diversity. Inspired by this gap, we demonstrate that steering vectors derived solely from text-only LLM backbones can effectively guide their multimodal counterparts, revealing a novel cross-modal transfer that enables reuse of existing interpretability tools. Using community-standard methods—Sparse Autoencoders (SAE), Mean Shift, and Linear Probing—we systematically validate this transfer effect across diverse MLLM architectures and visual reasoning tasks. Text-derived steering consistently enhances multimodal performance, with mean shift achieving up to +7.3% improvement in spatial relationship accuracy and +3.3% in counting accuracy on CV-Bench, and exhibits strong generalization to out-of-distribution datasets. These results highlight textual steering vectors as a powerful, efficient mechanism for enhancing grounding in MLLMs with minimal additional data collection and computational overhead.

1 INTRODUCTION

Steering large language models (LLMs) via their internal representations has emerged as a lightweight, interpretable paradigm for eliciting safe and controllable behavior (Li et al., 2023a; Turner et al., 2023; Sharkey et al., 2025, *inter alia*). However, similar steering approaches have not yet gained prominence for *multimodal large language models* (MLLMs). This is in part due to their relative recency, as well as the heterogeneity of their architectures compared to text-only LLMs. Moreover, many steering methods assume access to a dataset of contrast pairs (Marks and Tegmark, 2023) to construct steering vectors, which may not be readily available for multimodal inputs.

Our key finding is that internal representations from a text-only LLM backbone retain their steering effectiveness even after multimodal adaptation. This transfer effect enables a new multimodal steering paradigm that is agnostic to architecture and does not require specialized multimodal data. Importantly, it also allows us to directly repurpose steering methods originally developed for text-only models—such as Sparse Autoencoders (SAEs), Mean Shift, and Linear Probing—without modality-specific modifications. This bridges the mature ecosystem of text-based steering (McGrath et al., 2024; Durmus et al., 2024; Hanna et al., 2025) with the emerging space of multimodal models, providing a lightweight and interpretable pathway for enhancing multimodal reasoning.

Building on this insight, we propose a plug-and-play framework for multimodal steering. We extract steering vectors from text-only LLM backbones using established techniques and then apply them to the hidden states of their multimodal counterparts. This approach leverages the existing toolbox of steering methods, which have been extensively studied and evaluated in the text domain, to ensure accessibility and broader applicability for multimodal research. In contrast, developing new multimodal-specific steering methods would require both specialized datasets and bespoke implementations, which can be difficult to adapt across different modalities and fusion architectures.

We evaluate our approach across multiple open-weight MLLMs and a broad suite of visual reasoning tasks. Our method consistently outperforms prompting baselines—for example, Mean Shift achieves up to +7.3% improvement in spatial relationship accuracy on CV-Bench. Notably, while direct prompting is effective for controlling *text-only* LLM behavior (Wu et al., 2025), it provides little

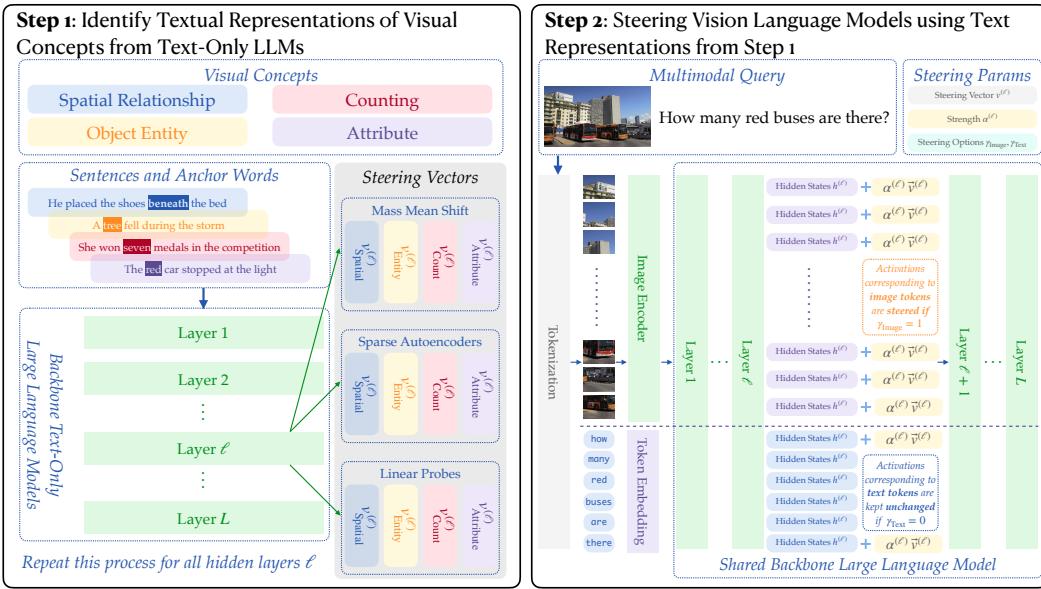


Figure 1: Overview of our steering methodology. Given an MLLM with a text-only LLM backbone and an image-bound prompt, we first identify the required visual concept (*e.g.*, spatial relationships, counting). For each hidden layer ℓ , we then extract corresponding steering vectors from the underlying LLM using Mean Shift, Linear Probing, or Sparse Autoencoders. Finally, we apply these vectors to image tokens, text tokens, or both, controlled by parameters γ_{Image} and γ_{Text} .

benefit for multimodal reasoning. We also compare against LoRA fine-tuning: although LoRA achieves stronger in-distribution accuracy, it exhibits limited out-of-distribution generalization and lacks the lightweight and interpretability advantages of steering. Our contributions are as follows:

- We introduce a plug-and-play multimodal steering method built directly on existing LLM representation-based techniques.
- We identify a novel transfer effect: representations from the text-only LLM backbone remain effective for steering its multimodal counterpart, even after vision-language post-training.
- We demonstrate consistent performance gains across multiple MLLMs and task categories. Importantly, we also show that textual steering vectors could generalize to out-of-distribution test sets and demonstrate significant performance gains.

2 RELATED WORKS

Representation-Based Steering methods are an effective family of methods for steering LLMs, often in two stages. First, they identify model components that influence target behaviors, using probing directions (Li et al., 2023a; Zou et al., 2023), activation differences (Li et al., 2023a; Turner et al., 2023; Panickssery et al., 2023; Marks and Tegmark, 2023; Lee et al., 2024), or lifted monosemantic features via SAEs (Lieberum et al., 2024b; Gao et al., 2025; Templeton et al., 2024; Marks et al., 2025) and their variants (Dunefsky et al., 2024), among other techniques. Second, they adjust steering hyperparameters to balance desiderata such as truthfulness (Lin et al., 2022; Hernandez et al., 2023; Li et al., 2023a), helpfulness (Zou et al., 2023), and quality.

While widely studied in LLMs, applying activation intervention to MLLMs remains elusive. To our knowledge, the only such effort is the VTI method (Liu et al., 2025), which extends LLM steering pipelines by constructing intervention vectors from paired multimodal inputs and applying them to both visual and textual representations. In contrast, we show that interventions vectors constructed solely from text inputs in the unimodal LLM can influence the MLLM’s multimodal behavior. This result highlights an underexplored form of cross-modal transfer enabled by the preserved semantics (Lieberum et al., 2024b) of the text backbone.

108 **Shared Semantics** refer to the representations unifying heterogeneous modalities of the same content,
 109 as identified across languages in multilingual LLMs (Artetxe et al., 2019; Wendler et al., 2024; Wu
 110 et al., 2024) and text/vision inputs in multimodal models (Huh et al., 2024; Luo et al., 2024; Wu
 111 et al., 2024). Our work studies the transfer of steering effect across different modalities and training
 112 stages. Concurrently, Papadimitriou et al. show that SAE features co-activate across multimodal
 113 inputs, while our work explores how such shared features can be exploited to steer MLLMs.

114 **Multimodal Large Language Models** are commonly developed by endowing a backbone LLM
 115 with visual processing components and fine-tuning on multimodal datasets, with some exceptions
 116 still pretrained from scratch (Team, 2024a; OLMo et al., 2024; Chen et al., 2025). Using an LLM
 117 backbone typically involves projecting the outputs of an image encoder (Dosovitskiy et al., 2020;
 118 Zhai et al., 2023) to the same dimension as the underlying LLM by an MLP, and concatenating the
 119 resulting image/text tokens as input to the LLM. The model can then be finetuned on multimodal
 120 data, possibly with frozen layers (*e.g.*, in the LLM) to preserve pretrained knowledge.

122 3 TOY EXAMPLE

124 To demonstrate that textual representations can effectively intervene in visual understanding, we conduct a
 125 simple color perception experiment using GemmaScope
 126 (Lieberum et al., 2024a) for Gemma-2-9B for feature ex-
 127 traction and PaliGemma2-10B-mix-448 (Beyer et al.,
 128 2024) as our target model. We present the model with a
 129 yellow-orange image (whose RGB hex code is `#FFB400`)
 130 and manipulate its perception by intervening in the hidden
 131 representations. Specifically, we obtain the normalized
 132 **red** vector from GemmaScope and we add this vector to
 133 the hidden states of image tokens at layer 20 as follows:
 134 $h'_{\text{image}} = h_{\text{image}} + \alpha \cdot v_{\text{red}}$, where α is the scale factor
 135 controlling intervention strength. Figure 2 shows how in-
 136 creasing the scale factor shifts perception along a color
 137 spectrum: initially yellow-orange dominates, then orange
 138 peaks at scale factor 50, and finally red becomes dominant
 139 beyond scale factor 75. This demonstrates that textual
 140 features can integrate with and modify visual understanding, supporting our hypothesis of unified
 141 cross-modal representations within these models. We include more color examples in Appendix B.

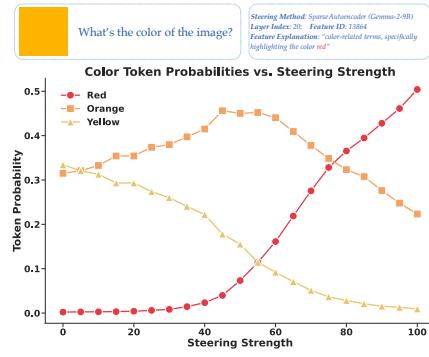
142 4 METHODS

144 Building on our demonstration that textual representations can effectively steer visual understanding,
 145 we now explore systematic approaches to improve MLLMs’ visual reasoning. Despite their growing
 146 success, MLLMs still struggle with seemingly simple visual queries—miscounting objects, confusing
 147 spatial relationships, and mishandling compositional prompts (Fu et al., 2024b). When the same
 148 problems are posed in pure text, foundation models perform far better (Fu et al., 2024a).

149 This observation motivates our central question: *Can existing steering mechanisms for textual*
 150 *representations rectify the shortcomings of MLLMs?* A promising remedy is steering vectors:
 151 compact directions in activation space that encode specific concepts. By adding these vectors to
 152 hidden representations at inference time as $h'_{\text{target}}^{(\ell)} = h_{\text{target}}^{(\ell)} + \alpha \cdot v^{(\ell)}$, we can amplify the model’s
 153 internal representation of desired concepts without parameter updates. The optimal layer ℓ^* and scale
 154 α^* are found via grid search. We use three established methods—Sparse Autoencoders (SAE), Mean
 155 Shift, and Linear Probing—to extract vectors $v^{(\ell)}$ from text-only LLM backbones, demonstrating
 156 broad applicability of cross-modal transfer while ensuring accessibility and reproducibility.

158 4.1 DATASET CONSTRUCTION FOR STEERING VECTOR EXTRACTION

160 To extract high-level textual representations for visual concepts, we identify four important tax-
 161 onomies for static images: spatial relationship, counting, attribute, and entity (Huang et al., 2023;
 Lin et al., 2024; Fu et al., 2025). For each visual concept, we curate small sets of sentence-anchor



162 Figure 2: Effect of steering strength on
 163 color token probabilities.

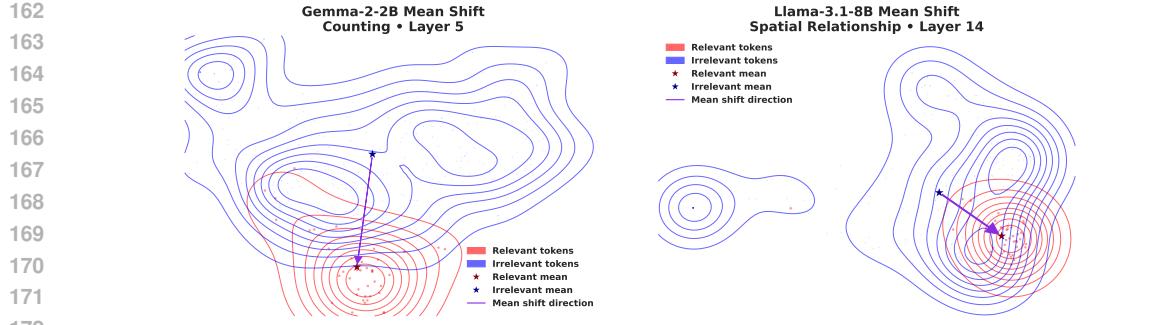


Figure 3: **Left:** Mean Shift method for counting features in Gemma-2-2B. The direction points from mean control token states to mean counting-related token states. **Right:** Spatial relationship features for Llama-3.1-8B. Activations projected to 2D for visualization.

pairs, where each pair contains a sentence exhibiting the visual concept and the specific anchor word representing that concept. These sentence-anchor pairs serve as the foundation for all three steering vector extraction methods. Examples are provided in Table 4 in Appendix A.1.

4.2 INTERPRETABLE STEERING VECTOR EXTRACTION METHODS

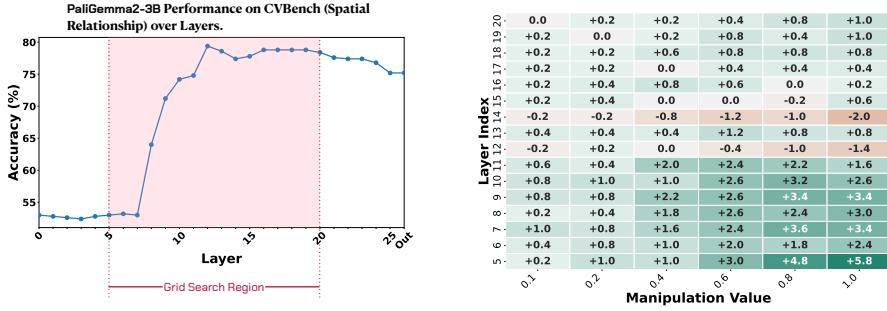
Sparse Autoencoders (SAE). Sparse Autoencoders reconstruct the activations of an LLM’s hidden layer using an MLP with a single hidden layer and a sparsity penalty on the hidden layer. More precisely, let $x = h^{(\ell)}(t) \in \mathbb{R}^D$ be the model activations for a token t at layer ℓ in an LLM. A SAE reconstructs x as $\hat{x} = b_{\text{dec}} + \sum_{i=1}^F f_i(x) W_{\cdot,i}^{\text{dec}}$, where $b_{\text{dec}} \in \mathbb{R}^D$ and $W^{\text{dec}} \in \mathbb{R}^{D \times F}$ are learned decoder weights, and $f_i(x)$ is the activation corresponding to feature i . Feature activations are computed using learned encoder weights $W^{\text{enc}} \in \mathbb{R}^{F \times D}$ and $b^{\text{enc}} \in \mathbb{R}^F$ as $f_i(x) = \sigma(W_{i,\cdot}^{\text{enc}} x + b_i^{\text{enc}})$, where σ denotes an activation function of choice, e.g., ReLU or JumpReLU.

The model is trained by minimizing the loss function $L = \mathbb{E}_x \left[\|x - \hat{x}\|_2^2 + \lambda \sum_{i=1}^F f_i(x) \|W_{\cdot,i}^{\text{dec}}\|_2 \right]$, i.e., L_2 -reconstruction error and L_1 -regularization on feature activations. In this formulation, unit-normalized decoder weight vectors $v_i^{(\ell)} := \frac{W_{\cdot,i}^{\text{dec}}}{\|W_{\cdot,i}^{\text{dec}}\|_2}$ serve as feature directions and $\alpha_i^{(\ell)}(t) := f_i(h^{(\ell)}(t)) \|W_{\cdot,i}^{\text{dec}}\|_2$ as the activation strength of $v_i^{(\ell)}$ on token t .

We leverage existing pretrained SAEs—GemmaScope (Lieberum et al., 2024b) for Gemma-2 models and LlamaScope (He et al., 2024) for Llama-3.1-8B. We emphasize that training SAEs is computationally expensive, and a key advantage of our approach is leveraging existing interpretability infrastructure without additional training costs. Using our sentence-anchor pairs, we identify features with high activations on anchor words. We then verify their relevance to the target visual concepts and average these relevant feature vectors to create a single steering vector for each visual concept at each layer. Additional details are provided in Appendix A.1.

Mean Shift. This method identifies feature directions by computing activation differences, as shown in Figure 3, showing surprising effectiveness for LLM steering (Marks and Tegmark, 2023; Wu et al., 2025). For each taxonomy \mathcal{T} and layer ℓ , using sentence-anchor pairs $\{(s_1, w_1), \dots, (s_K, w_K)\}$, we compute the mean shift vector $m_{\mathcal{T}}^{(\ell)} = \frac{1}{K} \sum_{j=1}^K h^{(\ell)}(w_j) - \frac{1}{|\mathcal{S}_{-\mathcal{T}}|} \sum_{t \in \mathcal{S}_{-\mathcal{T}}} h^{(\ell)}(t)$, where $h^{(\ell)}(w_j)$ represents the residual stream activation of the anchor word w_j at layer ℓ and $\mathcal{S}_{-\mathcal{T}}$ is a control set of non-anchor tokens from the same sentences. We refrain from normalizing the vector $m_{\mathcal{T}}^{(\ell)}$, preserving its magnitude relative to the original hidden states.

Linear Probing. We train a linear classifier distinguishing anchor word activations from control tokens on the ℓ -th layer of a model (Alain and Bengio, 2016; Park et al., 2024). As the hidden state dimensionality often exceeds our sample size ($K < D$), we first project to dimension $d < K$ using PCA. With $Q \in \mathbb{R}^{d \times D}$ as the PCA matrix, the probe separates $\{h^{(\ell)}(w_j) Q^\top\}_{j \leq K}$ and $\{h^{(\ell)}(t) Q^\top\}_{t \in \mathcal{S}_{-\mathcal{T}}}$, where $\{(s_1, w_1), \dots, (s_K, w_K)\}$ are the sentence-anchor pairs for concept \mathcal{T}



(a) We zero out models’ attention to image tokens after layer ℓ and measure model performance. This reveals when visual information is processed and allows efficient grid search.

(b) Grid search on PaliGemma2-3B to locate the best (ℓ^*, α^*) for steering the model’s spatial reasoning abilities. In this case, $\ell^* = 5$ and $\alpha^* = 1.0$.

Figure 4: Efficient Grid Search with PaliGemma2-3B on the Spatial Relationship Task.

and $\mathcal{S}_{-\mathcal{T}}$ is our control set. The learned normal vector $v \in \mathbb{R}^d$ (pointing toward taxonomy-relevant points) yields the final steering vector $v' = Q^\top v$. We use $d = K/2$ in practice.

Prompting Baseline. Like our steering methods, prompting represents an interpretable approach that no parameter updates, and it has displayed impressive steering abilities in text-only domains (Wu et al., 2025). For a given taxonomy \mathcal{T} , we generate a prompt meant to enhance an MLLM’s visual reasoning ability with respect to \mathcal{T} as follows: We first curate a collection of 96 prompts of varying lengths by instructing GPT-4o to generate prompts that guide the model to reason with respect to \mathcal{T} , similar to the LLM-based prompt generation in AxBench (Wu et al., 2025), and then select the best-performing prompt via grid search on training data. Refer to Appendix A.2 for further detail.

5 STEERING IMPROVES MULTIMODAL LLMs

Having established in Section 3 that textual steering vectors applied to non-output tokens can alter the behavior of MLLMs, we now investigate whether the textual steering vectors we identified in Section 4.2 can *improve* visual understanding in MLLMs when applied to intermediate representations.

5.1 SETUP

Models. We investigate PaliGemma2 models with 3B and 10B parameters (PaliGemma2-3B-mix-448 and PaliGemma2-10B-mix-448, referred to as PaliGemma2-3B and PaliGemma2-10B) and Idefics3-8B-Llama3. These models differ architecturally: PaliGemma2 adopts prefix-LM masking where image tokens and textual instructions are cross-attended, while Idefics3 is fully autoregressive following LLaVA. Steering vectors are extracted from their respective text-only backbones: Gemma2-2B, Gemma2-9B, and Llama-3.1-8B.

Dataset. We use CV-Bench (Tong et al., 2024) with 4 sub-categories: Count, Relation, Distance, and Depth, totaling 2,638 data points. Each sub-category contains around 700 samples, split into 500-600 training samples for grid search and 150 for testing.

Grid Search. We identify optimal injection layer ℓ and scale factor α via grid search on the training split. For each (ℓ, α) pair, we intervene as $h'_{\text{target}}(\ell) = h_{\text{target}}(\ell) + \alpha v^{(\ell)}$ and select $(\ell^*, \alpha^*) = \text{argmax}_{\ell, \alpha} \text{Acc}(\ell, \alpha)$. We use $\mathcal{A} = \{0.1, 0.2, 0.4, 0.6, 0.8, 1.0\}$ for unnormalized vectors (MeanShift). For normalized vectors (SAE, Probe), we use $\{10, 20, 30, 40, 50, 60\}$ on PaliGemma2 models and $\{0.2, 0.4, 0.6, 0.8, 1.0, 1.2\}$ on Idefics3 due to smaller hidden state norms. We set \mathcal{I} to be the middle layers, where we observe the learning from image tokens is predominantly happening (see Figure 4a): $\{5, 6, \dots, 20\}$ for PaliGemma2-3B and Idefics3-8B-Llama3, and $\{15, 16, \dots, 30\}$ for PaliGemma2-10B. Notably, we never steer output tokens, focusing on internal representations.

5.2 RESULTS

Table 1 presents a comparative analysis of three different models, PaliGemma2-3B, PaliGemma2-10B, and Idefics3-8B-Llama3, on tasks related to spatial relationships and counting in CV-Bench.

MODEL	INTERVENTION TOKENS		RELATION			COUNT		
	TEXT	IMAGE	SAE	PROBE	MEANSHIFT	SAE	PROBE	MEANSHIFT
PaliGemma2-3B	—	✓	82.0 (+6.0)*	77.3 (+1.3)	83.3 (+7.3)*	60.0 (+0.7)	62.0 (+2.7)	60.0 (+0.7)
	✓	✓	78.7 (+2.7)	76.7 (+0.7)	78.7 (+2.7)*	62.0 (+2.7)*	60.7 (+1.3)	62.0 (+2.7)
	✓	✓	81.3 (+5.3)*	78.7 (+2.7)	81.3 (+5.3)*	62.7 (+3.3)*	62.0 (+2.7)*	62.0 (+2.7)*
	Prompting	—	76.7 (+0.7)				60.0 (+0.7)	
PaliGemma2-10B	—	✓	79.3				63.3	
	✓	✓	78.7 (-0.7)	77.3 (-2.0)	83.3 (+4.0)*	63.3 (+0.0)	62.7 (-0.7)	64.0 (+0.7)
	✓	✓	79.3 (+0.0)	79.3 (+0.0)	78.7 (-0.7)	63.3 (+0.0)	63.3 (+0.0)	64.7 (+1.3)
	Prompting	—	78.7 (-0.7)	78.0 (-1.3)	83.3 (+4.0)*	64.0 (+0.7)	63.3 (+0.0)	63.3 (+0.0)
Idefics3-8B-Llama3	—	✓	73.3				59.3	
	✓	✓	76.0 (+2.7)	78.0 (+4.7)*	80.0 (+6.7)*	58.7 (-0.7)	58.0 (-1.3)	60.0 (+0.7)
	✓	✓	78.0 (+4.7)*	72.7 (-0.7)	76.7 (+3.3)	60.0 (+0.7)	59.3 (+0.0)	60.7 (+1.3)
	Prompting	—	77.3 (+4.0)*	78.7 (+5.3)*	80.7 (+7.3)*	62.0 (+2.7)*	60.0 (+0.7)	60.7 (+1.3)
			75.3 (+2.0)				58.7 (-0.7)	

Table 1: **Textual Steering Vectors Improve Multimodal LLMs’ Visual Understanding.** Task-specific textual steering vectors reliably improve both spatial relation and counting performance across models. Stars (*) denote statistically significant improvements ($p < 0.05$).

The performance is evaluated with and without intervention tokens (text, image, or both) and across different steering methods (SAE, Probe, MeanShift, and Prompting).

Steering Interventions Prove Effective. Table 1 demonstrates that steering interventions, especially MeanShift, consistently improve model performance on spatial relationship and counting tasks over baseline levels. For instance, PaliGemma2-3B’s “Relation” accuracy with MeanShift rose from 76.0 to 83.3 using both tokens, illustrating the general efficacy of these mechanisms.

MeanShift Shows Superior Performance and Stable Effects. Among the evaluated methods, MeanShift performs most effectively and demonstrates more stable effects across different models, aligning with recent text-only steering findings (Wu et al., 2025). MeanShift’s superiority and stability stem from its robustness: while SAE relies on learned sparse representations that may suffer from overfitting or incomplete concept capture, and probing operates in lower-dimensional space with sensitivity to specific projections, MeanShift operates on full-dimensional representations using distributional properties. This gives it more deterministic and stable effects across different models.

Prompting Barely Steers. Table 1 indicates that prompting is often less effective than targeted interventions and sometimes even deleterious. This deviates from text-only observations (Wu et al., 2025), reflecting MLLMs’ challenges in following fine-grained visual reasoning instructions. Unlike text-only models that reliably execute linguistic guidance, multimodal models may struggle with translating textual prompts into enhanced visual understanding, making prompting less effective.

Steering More Impactful for Spatial Relationships. Interventions yield more substantial accuracy improvements in the “Spatial Relationship” task than in “Counting”. For instance, as shown in Table 1, with both tokens and MeanShift, PaliGemma2-3B gained +7.3 for relationships but only +2.7 for counting. This disparity may stem from spatial relationships being more directly influenced by highlighting salient object features and positions, while counting might demand a more holistic scene interpretation, less directly aided by these specific steering methods.

Smaller Models Show Better Steering Responsiveness. Steering effectiveness increases with smaller model size, with the 3B model showing consistently larger improvements than the 10B model. This suggests that smaller models have more malleable internal representations, making them more receptive to steering interventions. For instance, PaliGemma2-3B demonstrates high responsiveness across all intervention types, while the 10B model shows reduced sensitivity to steering.

Intervention Transfers Across Tasks. As shown in Figure 5, intervention using a feature \mathcal{T} sometimes transfer effectively to different tasks \mathcal{T}' . For instance, enhancing attribute and entity recognition improves spatial relationship performance, suggesting that accurate object identification helps spatial reasoning. This cross-task transfer reflects the interconnected nature of visual understanding, where strengthening one capability can have cascading benefits for related reasoning processes.

	Counting	Spatial Relationship	Entity	Attribute	Counting	Spatial Relationship	Entity	Attribute	Counting	Spatial Relationship	Entity	Attribute		
Count	+0.7% (L7@0.8)	+1.3% (L14@1)	+0.0% (L9@0.8)	+0.0% (L16@0.6)	Count	+2.7% (L5@0.4)	+1.3% (L5@1)	+2.0% (L11@0.8)	+1.3% (L10@1)	Count	+2.7% (L9@0.6)	+1.3% (L5@0.6)	-0.7% (L7@0.2)	+2.0% (L10@1)
Relation	+3.3% (L9@1)	+7.3% (L5@1)	+0.0% (L9@0.6)	+2.7% (L5@1)	Relation	+0.7% (L10@0.8)	+2.7% (L10@1)	+3.3% (L14@1)	+4.0% (L13@1)	Relation	+3.3% (L9@1)	+5.3% (L5@1)	+0.7% (L6@0.8)	+2.0% (L6@0.8)
Distance	+1.3% (L6@0.8)	+0.7% (L16@0.6)	+0.7% (L6@0.1)	+0.7% (L5@0.2)	Distance	+0.0% (L19@0.8)	+0.0% (L9@0.6)	-0.7% (L6@0.2)	-2.0% (L10@0.4)	Distance	-0.7% (L15@0.8)	+2.0% (L8@0.2)	+1.3% (L6@0.1)	+1.3% (L13@0.2)
Depth	-0.7% (L11@0.6)	-1.3% (L11@0.4)	+0.7% (L10@0.8)	-0.7% (L10@0.4)	Depth	+0.7% (L11@1)	+2.7% (L10@1)	+0.7% (L11@0.8)	+0.7% (L10@1)	Depth	+0.0% (L5@0.8)	+2.0% (L11@0.6)	+1.3% (L10@0.8)	+0.0% (L10@0.8)

(a) Intervening Text Tokens

(b) Intervening Image Tokens

(c) Intervening Both Tokens

Figure 5: Performance improvements on CV-Bench tasks when steering PaliGemma2-3B with MeanShift vectors. Each cell shows the percentage improvement in accuracy relative to the baseline. Rows represent different CV-Bench tasks, while columns represent different feature vectors used for steering. Text below improvements indicates the optimal layer number and intervention strength.

6 STEERING IMPROVEMENTS GENERALIZE OUT-OF-DISTRIBUTION

We now examine the ability of textual steering methods for MLLMs to generalize out-of-distribution, *i.e.*, to datasets on which the steering method’s hyperparameters (ℓ, α) have not been tuned.

6.1 SETUP

Datasets. We first examine the transferability of textual steering on five datasets specifically designed to benchmark isolated visual reasoning capabilities: What’sUp-A, What’sUp-B, BLINK Object Localization, CLEVR, and Super-CLEVR. What’sUp-A contains 408 images of pairs of household objects arranged in clear spatial relations of {“on”, “under”, “left”, and “right”}, while What’sUp-B similarly contains 412 images with objects in the image closer in size (Kamath et al., 2023). The BLINK Object Localization category contains 122 questions related to bounding boxes for large objects (Fu et al., 2024b). Finally, we sampled 500 datapoints from CLEVR (Johnson et al., 2017) and 200 datapoints from Super-CLEVR (Li et al., 2023b) to evaluate the OOD accuracy of textual steering in counting.

Steering Vector Hyperparameter Selection. We examine the previous three steering methodologies—SAE, MeanShift, and Probe—with a single choice of layer ℓ and scale factor α chosen independently of the test dataset. Specifically, for each test dataset, we select the (ℓ, α) pair that performed best on the corresponding CV-Bench task category (*e.g.*, “Relation” for the What’sUp datasets and Blink Object Localization focusing on spatial relationships, and “Count” for CLEVR and Super-CLEVR).

We emphasize that the steering methods’ hyperparameters are *not* tuned to the datasets considered in this section, making this a true test of out-of-distribution generalization. Similarly, our prompting baseline uses the exact prompt prefix that performed best on the associated CV-Bench tasks. The only adaptation made was the use of a small validation subset (50 datapoints for What’sUp and CLEVR, 25 datapoints for BLINK Object Localization and Super-CLEVR) to determine the most effective token type for intervention (image, text, or both) before evaluating on the remaining data.

6.2 RESULTS

Steering Remains Broadly Effective. Table 2 demonstrates that textual interventions are effective across all 5 datasets considered, attaining an average improvement over all models and datasets of at least +3.9 for all vector-based steering methods, demonstrating the strong OOD generalization of steering. Prompting averaged a +0.8 improvement and worsened model performance in 5 cases, suggesting that it may be less effective for MLLMs than for text-only LLMs (Wu et al., 2025).

Validation Against Linguistic Bias. The improved performance of steering on the What’sUp datasets provides evidence that our steering enhances genuine visual understanding rather than exploiting linguistic patterns. These datasets contain controlled image groups where identical objects are arranged in systematically varied spatial relationships (*e.g.*, an apple positioned left, right, above, or

DATASET	VISUAL CONCEPT	MODEL	INTERVENTION METHOD				
			BASELINE	PROMPTING	SAE	PROBE	MEANSHIFT
What'sUp-A	Spatial Relation	PaliGemma2-3B	62.7	65.8 (+3.1)*	71.8 (+9.1)*	78.5 (+15.8)*	75.4 (+12.7)*
		PaliGemma2-10B	68.5	63.3 (-5.2)	80.1 (+11.6)*	71.6 (+3.1)*	74.9 (+6.4)*
		Idefics3-8B-Llama3	62.2	61.9 (-0.4)	64.1 (+1.9)	62.2 (+0.0)	61.9 (-0.3)
		AVERAGE IMPROVEMENT	-	-0.8	+7.6	+6.3	+6.3
What'sUp-B	Spatial Relation	PaliGemma2-3B	60.6	56.7 (-3.9)	58.9 (-1.7)	57.5 (-3.1)	60.3 (-0.3)
		PaliGemma2-10B	81.8	77.8 (-3.0)	82.4 (+0.6)	82.1 (+0.3)	82.1 (+0.3)
		Idefics3-8B-Llama3	52.0	57.3 (+5.3)*	56.2 (+4.2)*	57.0 (+5.0)*	63.4 (+11.5)*
		AVERAGE IMPROVEMENT	-	-0.5	+1.0	+0.8	+3.8
BLINK Object Localization	Spatial Relation	PaliGemma2-3B	41.2	41.2 (+0.0)	43.3 (+2.1)	42.3 (+1.0)	44.3 (+3.1)*
		PaliGemma2-10B	51.6	52.6 (+1.0)	54.6 (+3.1)	53.6 (+2.1)	57.7 (+6.2)*
		Idefics3-8B-Llama3	53.6	53.6 (+0.0)	56.7 (+3.1)*	53.6 (+0.0)	55.7 (+2.1)
		AVERAGE IMPROVEMENT	-	+0.3	+2.8	+1.0	+3.8
CLEVR	Count	PaliGemma2-3B	52.4	53.6 (+1.2)	70.7 (+18.2)*	56.4 (+4.0)*	67.1 (+14.7)*
		PaliGemma2-10B	70.7	72.4 (+1.7)	74.9 (+4.2)*	71.6 (+0.9)	80.4 (+9.8)*
		Idefics3-8B-Llama3	59.8	60.2 (+0.4)	88.0 (+28.2)*	84.4 (+24.7)*	94.0 (+34.2)*
		AVERAGE IMPROVEMENT	-	+1.1	+16.9	+9.9	+19.6
AVERAGE IMPROVEMENT			-	+0.8	+6.0	+3.9	+7.6

Table 2: Performance of textual steering on out-of-distribution datasets. Stars (*) denote statistically significant improvements ($p < 0.05$).

below the same plate). If we were merely exploiting textual patterns, we would expect biased outputs regardless of visual content, rather than the observed accurate tracking of true spatial relationships.

Superior OOD Performance on Focused Tasks. Remarkably, out-of-distribution performance often surpasses in-distribution results on CV-Bench, particularly on datasets requiring “pure” reasoning abilities. For example, CLEVR, which isolates counting skills using simple geometric objects without requiring complex object recognition, shows pronounced gains (+19.6 average). In contrast, CV-Bench Count and Super-CLEVR demand broader compositional understanding and object recognition beyond the targeted abilities, resulting in more moderate improvements. This pattern suggests our steering precisely targets the intended cognitive capabilities.

MeanShift Demonstrates Consistent Superiority. Across all experimental conditions, MeanShift consistently outperforms other extraction methods, achieving the highest average improvement of +7.6 compared to +6.0 for SAE and +3.9 for Linear Probing. This mirrors results from CV-Bench and AxBench (Wu et al., 2025), demonstrating MeanShift’s consistent superiority across different domains and modalities.

6.3 RESULTS ON REAL-WORLD TASKS

The datasets evaluated in the previous subsection were specifically designed to benchmark isolated visual reasoning capabilities—spatial relationships and counting—making them ideal for controlled evaluation of our steering methods. To examine broader practical applicability, we further evaluated our cross-modal steering approach on real-world multimodal tasks that MLLMs encounter in practical applications, including: general visual question answering (VQAv2 (Goyal et al., 2017)), open-ended image captioning (COCO Captions (Chen et al., 2015)), document understanding (DocVQA (Mathew et al., 2021)), chart understanding (ChartQA (Masry et al., 2022)), and table reasoning (VTabFact (Kim et al., 2024)). We applied the most conceptually related steering vectors to their corresponding tasks, using the same experimental protocol as Section 6.1. Both the experimental details and complete results are provided in Appendix C.

Our steering methods demonstrate consistent effectiveness across these diverse domains, with MeanShift achieving improvements in 15 out of 18 model-task combinations and 7 statistically significant gains. While improvement magnitudes are smaller than those observed on our primary OOD datasets, this is expected since these complex tasks depend less exclusively on core spatial relation and counting skills that our steering vectors specifically target. Despite these differences, the consistent positive

TASK	DATA TYPE	LoRA PERFORMANCE			AVERAGE IMPROVEMENT
		PALIGEMMA-3B	PALIGEMMA-10B	IDEFICS-8B	
CVBench Relation	In-dist	91.3 (+15.3)*	91.3 (+12.0)*	88.0 (+12.7)*	+13.3
CVBench Count	In-dist	67.3 (+8.0)*	72.0 (+8.7)*	67.3 (+8.0)*	+8.2
AVERAGE IN-DISTRIBUTION		+11.7	+10.4	+10.4	+10.8
What'sUp-A	OOD	67.7 (+5.0)*	69.3 (+0.8)	61.6 (-0.6)	+1.7
What'sUp-B	OOD	58.4 (-2.2)	86.0 (+4.2)*	58.1 (+6.1)*	+2.7
BLINK Object	OOD	42.3 (+1.1)	49.5 (-2.1)	52.6 (+1.0)	+0.0
CLEVR	OOD	54.2 (+1.8)	68.7 (-2.0)	66.7 (+6.9)*	+2.2
Super-CLEVR	OOD	28.6 (+1.7)	43.4 (+3.4)	66.9 (+0.4)	+1.8
AVERAGE OUT-OF-DISTRIBUTION		+1.3	+1.2	+2.8	+1.7

Table 3: Performance comparison between LoRA and baseline models across in-distribution and out-of-distribution tasks. Stars (*) denote statistically significant improvements ($p < 0.05$).

impact—especially the statistically significant gains—strongly indicates that our steering approach effectively enhances visual reasoning capabilities across diverse applications.

7 STEERING VS. FINE-TUNING

Beyond our interpretable steering methods, fine-tuning represents another common approach for enhancing model performance on specific tasks. To provide context for our steering approach, we compare against Low-Rank Adaptation (LoRA) fine-tuning (Hu et al., 2022) on the same tasks. We trained LoRA adapters using the training dataset from our grid search with an 80:20 train-validation split with hyperparameters: rank $r \in \{1, 2, 4\}$, alpha $\alpha \in \{4, 8\}$, learning rate $\eta \in \{1 \times 10^{-5}, 5 \times 10^{-5}, 1 \times 10^{-4}\}$, epochs = 3, and dropout = 0.1. We applied LoRA to the query and value projection parameters at the same layers used in our grid search: layers 5-20 for PaliGemma2-3B and Idefics3-8B, and layers 15-30 for PaliGemma2-10B. For each model and task combination, we selected hyperparameters that achieved optimal validation performance.

Table 3 presents the performance comparison between LoRA fine-tuning and our baseline models across in-distribution and out-of-distribution tasks. LoRA demonstrates strong in-distribution performance with an average improvement of +10.8 on CV-Bench, but its effectiveness diminishes significantly on out-of-distribution datasets with only +1.7 average improvement. In contrast, our steering methods maintain consistent performance across diverse datasets, with MeanShift achieving +7.6 and SAE achieving +6.0 average out-of-distribution improvements, highlighting the superior generalization capabilities of steering. This performance differential reflects fundamental differences in their mechanisms: LoRA adapts models to specific task distributions, while steering enhances underlying cognitive abilities such as spatial reasoning that remain applicable across diverse contexts.

8 DISCUSSION

We examine the ability of multimodal large language models (MLLMs) to be steered using textual steering vectors from their text-only backbone. We find that vectors extracted from Sparse Autoencoders (SAEs), Mean Shift, and Linear Probing can all enhance MLLMs’ visual reasoning across diverse tasks on CV-Bench, with Mean Shift demonstrating the strongest overall performance. Notably, steering vectors with hyperparameters optimized on CV-Bench generalize to other out-of-distribution datasets with superior performance compared to LoRA fine-tuning or prompt tuning, underscoring text-driven steering as a powerful and efficient medium for enhancing visual reasoning in MLLMs. A primary limitation of our steering method is the reliance on the quality of extracted steering vectors. While existing vector extraction methods are widely used in the LLM interpretability community, the vectors they extract can be of poor quality and fail to adequately represent the target concepts, particularly for SAE and Linear Probing, leading to variable steering performance across different layers and models. Future work can focus on developing more robust extraction methods for text-only or cross-modal models to improve the reliability and consistency of steering vectors.

486 **ETHICS STATEMENT**
487

488 We identify no significant ethical concerns. Our steering methods enhance visual reasoning on
489 standard benchmarks without introducing inherently harmful capabilities. While these techniques
490 could potentially be misused like any model modification approach, the risk is not greater than that of
491 the underlying MLLMs. We encourage responsible use and consideration of societal impacts when
492 deploying enhanced MLLMs.

493
494 **REPRODUCIBILITY STATEMENT**
495

496 We provide comprehensive implementation details and have open-sourced our code on GitHub and
497 uploaded it as supplementary material to OpenReview. Steering vector extraction methods are detailed
498 in Section 4.2 and Algorithm 1. Hyperparameter grid search procedures and experimental protocols
499 are described in Sections 5.1, 5, and 6. We use publicly available pre-trained SAEs (GemmaScope,
500 LlamaScope), models (PaliGemma2, Idefics3), and datasets (CV-Bench, What'sUp, BLINK, CLEVR,
501 Super-CLEVR, VQAv2, COCO Captions, DocVQA, ChartQA, VTabFact).

502
503 **REFERENCES**
504

505 Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes.
506 *arXiv preprint arXiv:1610.01644*, 2016.

507
508 Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. On the cross-lingual transferability of monolin-
509 gual representations. *arXiv preprint arXiv:1910.11856*, 2019.

510
511 Llama Team at Meta. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.

512
513 Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel
514 Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschanen, Emanuele Bugliarello, et al.
515 Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024.

516
517 Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and
518 Chong Ruan. Janus-pro: Unified multimodal understanding and generation with data and model
519 scaling. *arXiv preprint arXiv:2501.17811*, 2025.

520
521 Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and
522 C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint*
523 *arXiv:1504.00325*, 2015.

524
525 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
526 Unterthiner, Mostafa Dehghani, Matthias Minderer, G Heigold, S Gelly, et al. An image is
527 worth 16x16 words: Transformers for image recognition at scale. In *International Conference on*
528 *Learning Representations*, 2020.

529
530 Jacob Dunefsky, Philippe Chlenski, and Neel Nanda. Transcoders find interpretable llm feature
531 circuits. *arXiv preprint arXiv:2406.11944*, 2024.

532
533 Esin Durmus, Alex Tamkin, Jack Clark, Jerry Wei, Jonathan Marcus, Joshua Batson, Kunal
534 Handa, Liane Lovitt, Meg Tong, Miles McCain, Oliver Rausch, Saffron Huang, Sam Bow-
535 man, Stuart Ritchie, Tom Henighan, and Deep Ganguli. Evaluating feature steering: A
536 case study in mitigating social biases, 2024. URL <https://anthropic.com/research/evaluating-feature-steering>.

537
538 Deqing Fu, Ruohao Guo, Ghazal Khalighinejad, Ollie Liu, Bhuwan Dhingra, Dani Yogatama,
539 Robin Jia, and Willie Neiswanger. Isobench: Benchmarking multimodal foundation models on
540 isomorphic representations. In *First Conference on Language Modeling*, 2024a. URL <https://openreview.net/forum?id=KZd1EErJ1>.

-
- 540 Deqing Fu, Tong Xiao, Rui Wang, Wang Zhu, Pengchuan Zhang, Guan Pang, Robin Jia, and
541 Lawrence Chen. TLDR: Token-level detective reward model for large vision language models.
542 In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=Zy2XgaGpDw>.
543
- 544 Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith,
545 Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not
546 perceive. In *European Conference on Computer Vision*, pages 148–166. Springer, 2024b.
547
- 548 Leo Gao, Tom Dupre la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever,
549 Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. In *The Thirteenth
550 International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=tcsZt9ZNKD>.
551
- 552 Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa
553 matter: Elevating the role of image understanding in visual question answering. In *Proceedings of
554 the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017.
555
- 556 Michael Hanna, Mateusz Piotrowski, Jack Lindsey, and Emmanuel Ameisen. circuit-tracer. <https://github.com/safety-research/circuit-tracer>, 2025. The first two authors contributed
557 equally and are listed alphabetically.
558
- 559 Zhengfu He, Wentao Shu, Xuyang Ge, Lingjie Chen, Junxuan Wang, Yunhua Zhou, Frances Liu,
560 Qipeng Guo, Xuanjing Huang, Zuxuan Wu, et al. Llama scope: Extracting millions of features
561 from llama-3.1-8b with sparse autoencoders. *arXiv preprint arXiv:2410.20526*, 2024.
562
- 563 Evan Hernandez, Belinda Z Li, and Jacob Andreas. Inspecting and editing knowledge representations
564 in language models. *arXiv preprint arXiv:2304.00740*, 2023.
565
- 566 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
567 Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
568
- 569 Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2i-compbench: A com-
570 prehensive benchmark for open-world compositional text-to-image generation. *Advances in Neural
571 Information Processing Systems*, 36:78723–78747, 2023.
572
- 573 Minyoung Huh, Brian Cheung, Tongzhou Wang, and Phillip Isola. The platonic representation
574 hypothesis. *arXiv preprint arXiv:2405.07987*, 2024.
575
- 576 Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick, and
577 Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual
578 reasoning. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE,
579 2017.
580
- 581 Amita Kamath, Jack Hessel, and Kai-Wei Chang. What’s “up” with vision-language models?
582 investigating their struggle with spatial reasoning. In *Proceedings of the 2023 Conference on
583 Empirical Methods in Natural Language Processing*, pages 9161–9175, 2023.
584
- 585 Yoonsik Kim, Moonbin Yim, and Ka Yeon Song. Tablevqa-bench: A visual question answering
586 benchmark on multiple table domains. *arXiv preprint arXiv:2404.19205*, 2024.
587
- 588 Bruce W Lee, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Erik Miehling, Pierre Dognin, Manish
589 Nagireddy, and Amit Dhurandhar. Programming refusal with conditional activation steering. *arXiv
590 preprint arXiv:2409.05907*, 2024.
- 591 Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time
592 intervention: Eliciting truthful answers from a language model. In *Thirty-seventh Conference on
593 Neural Information Processing Systems*, 2023a. URL <https://openreview.net/forum?id=aLLuYpn83y>.
594
- 595 Zhuowan Li, Xingrui Wang, Elias Stengel-Eskin, Adam Kortylewski, Wufei Ma, Benjamin
596 Van Durme, and Alan L. Yuille. Super-clevr: A virtual benchmark to diagnose domain ro-
597 bustness in visual reasoning. In *Proceedings of the IEEE/CVF Conference on Computer Vision
598 and Pattern Recognition (CVPR)*, pages 14963–14973, June 2023b.
599

- 594 Tom Lieberum, Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant
595 Varma, János Kramár, Anca Dragan, Rohin Shah, and Neel Nanda. Gemma scope: Open sparse
596 autoencoders everywhere all at once on gemma 2. *arXiv preprint arXiv:2408.05147*, 2024a.
597
- 598 Tom Lieberum, Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant
599 Varma, János Kramár, Anca Dragan, Rohin Shah, and Neel Nanda. Gemma scope: Open sparse
600 autoencoders everywhere all at once on gemma 2, 2024b. URL <https://arxiv.org/abs/2408.05147>.
601
- 602 Stephanie Lin, Jacob Hilton, and Owain Evans. TruthfulQA: Measuring how models mimic human
603 falsehoods. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of
604 the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,
605 pages 3214–3252, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi:
606 10.18653/v1/2022.acl-long.229. URL <https://aclanthology.org/2022.acl-long.229/>.
607
- 608 Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and
609 Deva Ramanan. Evaluating text-to-visual generation with image-to-text generation. In *European
610 Conference on Computer Vision*, pages 366–384. Springer, 2024.
- 611 Sheng Liu, Haotian Ye, and James Zou. Reducing hallucinations in large vision-language models via
612 latent space steering. In *The Thirteenth International Conference on Learning Representations*,
613 2025. URL <https://openreview.net/forum?id=LB17Hez0ff>.
- 614 Grace Luo, Trevor Darrell, and Amir Bar. Task vectors are cross-modal. *arXiv preprint
615 arXiv:2410.22330*, 2024.
- 616 Samuel Marks and Max Tegmark. The geometry of truth: Emergent linear structure in large language
617 model representations of true/false datasets. *arXiv preprint arXiv:2310.06824*, 2023.
- 618 Samuel Marks, Can Rager, Eric J Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller.
619 Sparse feature circuits: Discovering and editing interpretable causal graphs in language models.
620 In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=I4e82CIDxv>.
621
- 622 Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A bench-
623 mark for question answering about charts with visual and logical reasoning. *arXiv preprint
624 arXiv:2203.10244*, 2022.
- 625 Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document
626 images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*,
627 pages 2200–2209, 2021.
- 628 Thomas McGrath, Daniel Balsam, Myra Deng, and Eric Ho. Understanding and steering llama 3
629 with sparse autoencoders. *Goodfire Research*, September 2024.
- 630 Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia,
631 Yuling Gu, Shengyi Huang, Matt Jordan, et al. 2 olmo 2 furious. *arXiv preprint arXiv:2501.00656*,
632 2024.
- 633 OpenAI. o3-mini. <https://openai.com/index/openai-o3-mini/>, 2025. Accessed: 2025-05-
634 13.
- 635 Nina Panickssery, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt
636 Turner. Steering llama 2 via contrastive activation addition. *arXiv preprint arXiv:2312.06681*,
637 2023.
- 638 Isabel Papadimitriou, Huangyuan Su, Thomas Fel, Sham Kakade, and Stephanie Gil. Interpreting the
639 linear structure of vision-language model embedding spaces. *arXiv preprint arXiv:2504.11695*,
640 2025.
- 641 Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry
642 of large language models. In *International Conference on Machine Learning*, pages 39643–39666.
643 PMLR, 2024.

-
- 648 Lee Sharkey, Bilal Chughtai, Joshua Batson, Jack Lindsey, Jeff Wu, Lucius Bushnaq, Nicholas
649 Goldowsky-Dill, Stefan Heimersheim, Alejandro Ortega, Joseph Bloom, et al. Open problems in
650 mechanistic interpretability. *arXiv preprint arXiv:2501.16496*, 2025.
- 651
- 652 Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint*
653 *arXiv:2405.09818*, 2024a.
- 654 Gemma Team. Gemma 2: Improving open language models at a practical size, 2024b. URL
655 <https://arxiv.org/abs/2408.00118>.
- 656
- 657 Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam
658 Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner,
659 Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees,
660 Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity:
661 Extracting interpretable features from claude 3 sonnet. *Transformer Circuits Thread*,
662 2024. URL <https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html>.
- 663
- 664 Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula,
665 Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, Ziteng Wang, Rob Fergus, Yann LeCun,
666 and Saining Xie. Cambrian-1: A fully open, vision-centric exploration of multimodal llms, 2024.
667 URL <https://arxiv.org/abs/2406.16860>.
- 668
- 669 Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J Vazquez, Ulisse Mini,
670 and Monte MacDiarmid. Steering language models with activation engineering. *arXiv preprint*
arXiv:2308.10248, 2023.
- 671
- 672 Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. Do llamas work in english?
673 on the latent language of multilingual transformers. In *Proceedings of the 62nd Annual Meeting*
674 *of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15366–15394,
675 2024.
- 676
- 677 Zhaofeng Wu, Xinyan Velocity Yu, Dani Yogatama, Jiasen Lu, and Yoon Kim. The semantic hub
678 hypothesis: Language models share semantic representations across languages and modalities.
arXiv preprint arXiv:2411.04986, 2024.
- 679
- 680 Zhengxuan Wu, Aryaman Arora, Atticus Geiger, Zheng Wang, Jing Huang, Dan Jurafsky, Christo-
681 pher D Manning, and Christopher Potts. Axbench: Steering llms? even simple baselines outperform
682 sparse autoencoders. *arXiv preprint arXiv:2501.17148*, 2025.
- 683
- 684 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language
685 image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*,
686 pages 11975–11986, 2023.
- 687
- 688 Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan,
689 Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, et al. Representation engineering: A
690 top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*, 2023.
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701

APPENDIX

702	A Steering Vector Methodology	14
703	A.1 Sparse Autoencoders	14
704	A.2 Prompting	15
705	B Additional Color Perception Intervention Examples	18
706		
707	C Results on Real-World Tasks	19
708		
709	D Dataset Evaluation Details	20
710		
711	E LLM Usage Statement	28
712		

A STEERING VECTOR METHODOLOGY

A.1 SPARSE AUTOENCODERS

We now provide further detail regarding the extraction of textual steering vectors for visual concepts using SAEs.

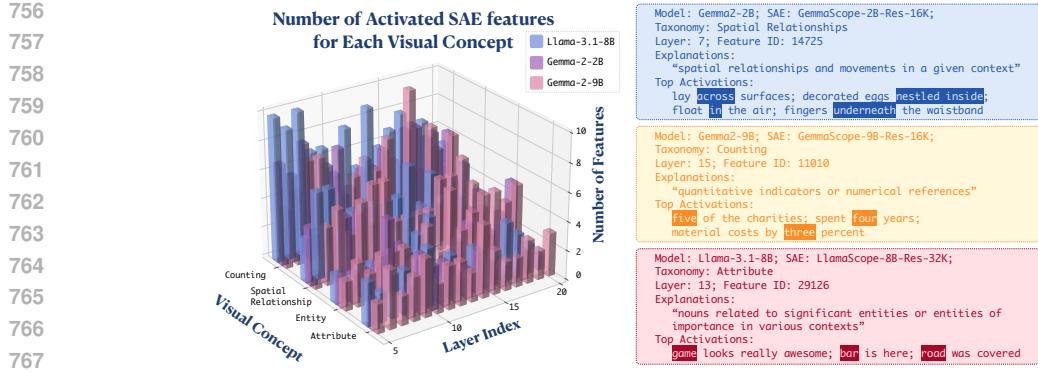
Recall that we consider four important taxonomies for image-related concepts: spatial relationship, counting, attribute, and entity. For each taxonomy, we sample K sentences $\{s_1, \dots, s_K\}$ containing these visual concepts. In practice, we set K to 20. For each sentence s_j , we identify the anchor word for this visual concept as w_j , thus forming sentence-anchor pairs (s_j, w_j) . See table 4 for several examples.

Table 4: Sample sentence and anchor word pairs for various taxonomies.

TAXONOMY	SENTENCE s_j	ANCHOR WORD w_j
Spatial Relationship	The cat is on the table	on
	She put the book under the chair	under
Counting	There are three apples in the basket	three
	The teacher counted five children	five
Attribute	The red car stopped at the light	red
	She wore a beautiful dress	beautiful
Entity	The dog barked at the mailman	dog
	A tree fell during the storm	tree

Using our sentence-anchor pairs, we identify features with high activations on anchor words. Interestingly, as shown in Figure 6, we find that each visual concept activates only a limited number of SAE features, indicating a sparse encoding of these concepts. We then verify their relevance to the target visual concepts and average these relevant feature vectors to create a single steering vector for each visual concept at each layer.

We then use these sentence-anchor pairs to identify feature directions corresponding to the ideal visual concepts using Algorithm 1. We employ a two-stage procedure which, at the first stage, finds the top n activated features for anchor words w_j in sentences s_j . At the second stage, we use o3-mini (OpenAI, 2025) to verify that these features indeed align with the desired visual concept \mathcal{C} . To accomplish the procedure, we use pretrained SAEs with detailed explanations and top activations developed by the interpretability community, such as GemmaScope (Lieberum et al., 2024b) for Gemma-2-2B and Gemma-2-9B (Team, 2024b), and LlamaScope (He et al., 2024) for Llama-3.1-8B base model (at Meta, 2024). When we prompt o3-mini for verification, we craft



756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

Figure 6: **Left:** Number of SAE features associated with each taxonomy (counting, spatial relationship, entity, and attribute) across the layers of Llama-3.1-8B, Gemma2-2B, and Gemma2-9B. Notably, SAE features for such visual concepts are sparse, numbering fewer than 10 across 16k total SAE features (Gemma2-2B/9B) or 32k features (Llama-3.1-8B). **Right:** Examples of features corresponding to visual concepts, identified by the layer whose activation space they inhabit and their (arbitrary) feature ID. The feature’s explanation summarizes its semantic meaning, as evidenced by the tokens and contexts on which it attains the greatest activations.

prompts to include both the explanation for the candidate feature vector $v_i^{(\ell)}$, and sample top activated tokens (see figure 7 for the prompting template). We find that o3-mini can indeed filter out features unrelated to the desired visual concepts.

Algorithm 1 Find Textual Representations for Visual Concepts using SAEs

Require: Desired visual concepts \mathcal{C} . Layer index ℓ .
Require: Sentence and anchor word pairs $\{(s_1, w_1), \dots, (s_K, w_K)\}$.
Require: Pretrained SAEs at layer ℓ .

▷ Find top activations and their corresponding SAE feature vectors.
 $\mathcal{V}_0 = \{\}$
for each (s_j, w_j) **do**
 $\{\alpha_i^{(\ell)}(w_j), v_i^{(\ell)}\} \leftarrow$ Pass s_j into the pretrained SAE
 $\{v_{i_1}^{(\ell)}, \dots, v_{i_n}^{(\ell)}\} \leftarrow$ Top n $\{\alpha_i^{(\ell)}(w_j), v_i^{(\ell)}\}$ ranked by activation strength $\alpha_i^{(\ell)}(w_j)$
 $\mathcal{V}_0 \leftarrow \mathcal{V}_0 \cup \{v_{i_1}^{(\ell)}, \dots, v_{i_n}^{(\ell)}\}$
end for
▷ Filter out noisy SAE feature vectors.
 $\mathcal{V} = \{\}$
for each $v_i^{(\ell)} \in \mathcal{V}_0$ **do**
 Find the explanation e and top activated tokens $\{t_1, \dots, t_p\}$ for $v_i^{(\ell)}$
 if o3-mini(VerificationPrompt, e , $\{t_1, \dots, t_p\}$, \mathcal{C}) is True **then**
 $\mathcal{V} \leftarrow \mathcal{V} \cup \{v_i^{(\ell)}\}$
 end if
end for
▷ Aggregate SAE vectors to one steering vector.
 $v^{(\ell)} = \frac{1}{|\mathcal{V}|} \sum_{u \in \mathcal{V}} u$
return $v^{(\ell)}$

806
807
A.2 PROMPTING

We now elaborate upon our generation of prompts for eliciting taxonomy-specific visual reasoning in MLLMs. As described in Section 4.2, we generate a total of 96 candidate prompts for each taxonomy \mathcal{T} . To do so, we use template shown in figure 8. Here, we set the num instructions to 6 and word

810 **FEATURE ALIGNMENT VERIFICATION**

811

812 Task: Determine if a neural network's sparse autoencoder (SAE)

813 feature aligns with the taxonomy "{taxonomy}".

814

815 Taxonomy Definition: {taxonomy_definition}

816

817 Feature Information:

818 1. Feature's explanation: {feature_explanation}

819 2. Top activation examples (tokens wrapped in <top>...</top> have the

820 highest activation values and are the most important to focus on):

821 1. {activation_example_1}

822 2. {activation_example_2}

823 3. {activation_example_3}

824 4. {activation_example_4}

825 5. {activation_example_5}

826

827 Examples of features that DO align with the {taxonomy} taxonomy

828 (Notice how the key words are highlighted with <top>...</top> tags):

829 Example 1:

830 - Explanation: {explanation_1}

831 - Activations: {activations_1}

832 Example 2:

833 - Explanation: {explanation_2}

834 - Activations: {activations_2}

835

836 When making your decision, you should follow these rules:

837 1. First pay attention to the feature's explanation.

838 2. If you cannot decide, you should then pay special attention to

839 the tokens highlighted with <top>...</top> tags, as these are the most

840 highly activated tokens and strongest indicators of what the feature

841 detects.

842 3. Also consider the diversity of the activation examples provided.

843 If one feature only activates one particular word, it may not be as

844 aligned as a feature that activates on a variety of words.

845

846 Based on the feature's explanation and the highlighted tokens in the

847 activation examples, does this feature specifically detect or respond

848 to {taxonomy_definition}? Your answer should start with YES or NO,

849 then provide a brief reason. Do not start with any other words or

850 phrases such as 'answer'.

849
850 Figure 7: Prompt template for querying GPT-o3-mini to verify whether a given feature is related to
851 a visual taxonomy. For each taxonomy, the template employs a brief definition of the taxonomy, two
852 example features that align with each taxonomy (for few-shot learning), and the top five activations
853 of the feature in question.

854
855 count ∈ {5, 10, 15, 20, 25, 30, 35, 40, 45, 50, 55, 60, 65, 70, 75, 80}, resulting in total 6 × 16 = 96
856 steering prompts.

857
858
859
860
861
862
863

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878

STEERING PROMPT GENERATION

System prompt: You are an expert at creating concise, clear instructions for Multimodal Large Language Models (MLLM).

Your task:

- Generate {num_instructions} different instruction(5) that will make the Model focus on {concept} when answering questions about images
 - Each instruction must be within {word_count} words
 - Instructions should be direct and actionable, focusing specifically on how to emphasize {concept}

IMPORTANT FORMAT REQUIREMENTS:

- Begin each instruction with "INSTRUCTION:" followed by the instruction text
 - Put each instruction on its own line
 - Do not include any numbering, bullets, or other text beyond the requested instructions
 - Do not include any explanations, introductions, or conclusions

Example format for 2 instructions:

INSTRUCTION: First instruction text here within word limit.

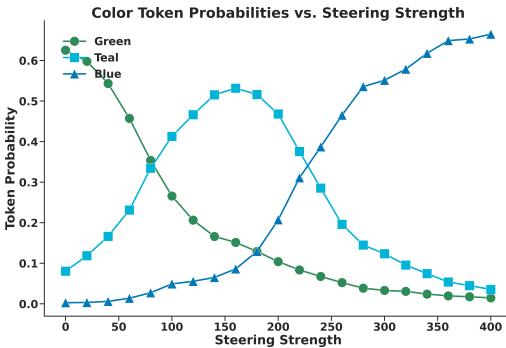
INSTRUCTION: Second instruction text here within word limit.

User prompt: Create {num_instructions} instruction(s) about {concept} using {word_count} words or fewer each.

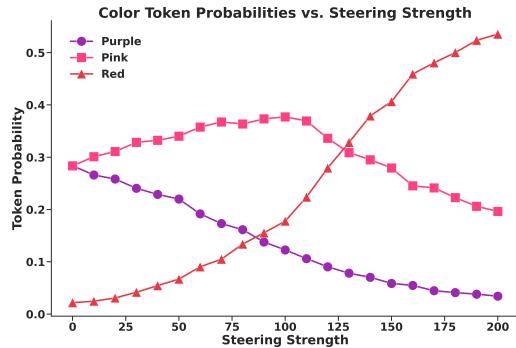
Figure 8: System and user prompt template for generating MLLM prompts.

B ADDITIONAL COLOR PERCEPTION INTERVENTION EXAMPLES

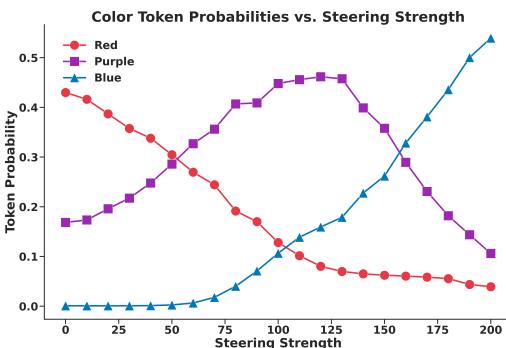
To further demonstrate the effectiveness of textual steering vectors in modifying visual understanding within MLLMs, we present additional color perception intervention examples using the same methodology described in §3.



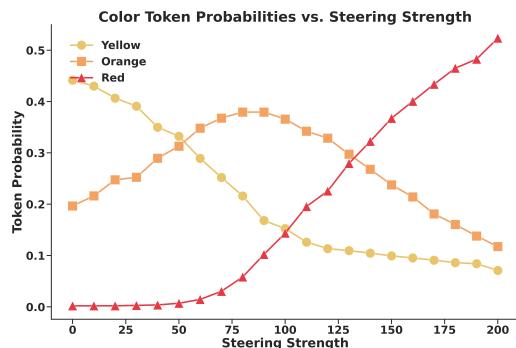
(a) Steering a green image toward blue perception. As the scale factor increases, the model’s interpretation shifts from green to teal, and ultimately to blue.



(b) Steering a purple image toward red perception. The intervention gradually shifts the model’s color association from purple to pink, and finally to red.



(c) Steering an red image toward blue perception. The intervention causes a gradual shift from red to purple, and ultimately to blue.



(d) Another example of steering a yellow image toward red perception, using a different steering vector from layer 18 of PaliGemma2-10B. As the scale factor increases, the model’s interpretation transitions from yellow to orange, and finally to red.

Figure 9: Additional color perception intervention examples. In each case, we apply the normalized textual steering vector for the target color to the image tokens with increasing scale factors. The steering vectors are extracted from and applies to one selected layer from layer 17 to 20 in PaliGemma2-10B. The plots show token probability shifts, demonstrating how textual steering vectors can systematically modify the model’s visual perception.

These additional examples further support our findings in §3. In each case, we see a clear progression of perception as the steering strength increases, with intermediate colors appearing during the transition. This confirms that textual steering vectors can produce predictable and continuous modifications to visual understanding.

Notably, all these interventions were performed using steering vectors derived solely from text data, yet they effectively modulate multimodal understanding. This provides additional evidence for our hypothesis that MLLMs develop unified cross-modal representations that can be manipulated through textual steering.

972 973 974 975 976 977 978 979 980 981 982 983 984 985 986 987 988 989 990	TASK	VISUAL CONCEPT	MODEL	INTERVENTION METHOD				
				BASELINE	PROMPTING	SAE	PROBE	MEANSHIFT
VQAv2	Spatial Relations	PaliGemma2-3B	86.8	87.0 (+0.2)	88.2 (+2.4)	87.1 (+0.3)	89.3 (+3.5)*	
		PaliGemma2-10B	88.2	86.8 (-1.4)	87.4 (-0.8)	86.9 (-1.3)	88.9 (+0.7)	
		Idefics3-8B	76.7	76.7 (+0.0)	78.1 (+1.4)	77.8 (+1.1)	74.6 (-2.1)	
		AVERAGE IMPROVEMENT	-	-0.4	+1.0	+0.0	+0.7	
COCO Captions	Spatial Relations	PaliGemma2-3B	147.9	144.4 (-3.5)	151.2 (+3.3)*	151.0 (+3.1)*	152.5 (+4.6)*	
		PaliGemma2-10B	155.8	141.3 (-14.5)	160.0 (+4.2)*	161.1 (+5.3)*	158.4 (+2.6)	
		Idefics3-8B	70.0	70.3 (+0.3)	71.2 (+1.2)	70.9 (+0.9)	69.6 (-0.5)	
		AVERAGE IMPROVEMENT	-	-5.9	+2.9	+3.1	+2.2	
DocVQA Layout	Spatial Relations	PaliGemma2-3B	79.4	81.4 (+2.0)	84.8 (+5.4)*	81.0 (+1.6)	85.4 (+6.0)*	
		PaliGemma2-10B	81.3	82.5 (+1.2)	83.8 (+2.5)	83.9 (+2.6)*	83.8 (+2.5)*	
		Idefics3-8B	88.5	86.3 (-2.2)	89.6 (+1.1)	88.2 (-0.3)	89.7 (+1.3)	
		AVERAGE IMPROVEMENT	-	+0.3	+3.0	+1.3	+3.3	
DocVQA Number	Counting	PaliGemma2-3B	76.1	76.2 (+0.1)	75.8 (-0.3)	76.3 (+0.2)	76.5 (+0.4)	
		PaliGemma2-10B	77.7	75.8 (-1.9)	77.6 (-0.1)	79.4 (+1.7)	76.9 (-0.8)	
		Idefics3-8B	86.8	84.5 (-2.3)	87.3 (+0.5)	86.9 (+0.1)	89.8 (+3.0)	
		AVERAGE IMPROVEMENT	-	-1.4	+0.0	+0.7	+0.9	
ChartQA	Counting	PaliGemma2-3B	46.4	45.4 (-1.0)	46.6 (+0.2)	47.4 (+1.0)	48.0 (+1.6)	
		PaliGemma2-10B	51.8	53.2 (+1.4)	53.8 (+2.0)	53.4 (+1.6)	54.4 (+2.6)*	
		Idefics3-8B	68.2	67.4 (-0.8)	71.0 (+2.8)*	67.2 (-1.0)	72.6 (+4.4)*	
		AVERAGE IMPROVEMENT	-	-0.1	+1.7	+0.5	+2.9	
VTabFact	Counting	PaliGemma2-3B	56.5	54.5 (-2.0)	58.0 (+1.5)	56.0 (-0.5)	60.5 (+4.0)*	
		PaliGemma2-10B	57.0	58.5 (+1.5)	58.5 (+1.5)	59.0 (+2.0)	58.5 (+1.5)	
		Idefics3-8B	70.0	71.0 (+1.0)	75.5 (+5.5)*	71.0 (+1.0)	73.5 (+3.5)	
		AVERAGE IMPROVEMENT	-	+0.2	+2.8	+0.8	+3.0	

991
992 Table 5: Performance of textual steering methods on real-world multimodal tasks. Stars (*) denote
993 statistically significant improvements ($p < 0.05$).

994 995 C RESULTS ON REAL-WORLD TASKS

996
997 **Experimental Setup:** We evaluated our steering methods on six real-world multimodal tasks using
998 500 examples per dataset (200 for VTabFact due to dataset size limitations) for testing, and extra 50
999 examples for validation to determine optimal intervention token types. We applied counting steering
1000 vectors to numerical reasoning tasks and spatial relationship vectors to layout and captioning tasks.

1001 Task Details and Metrics:

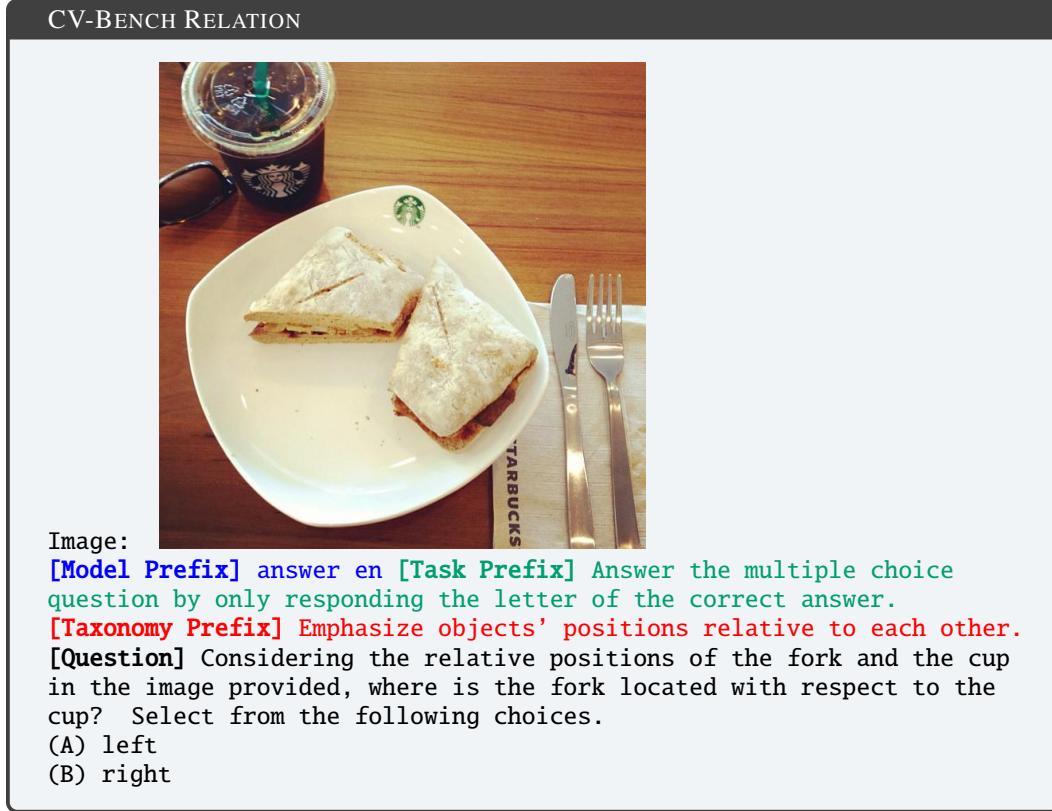
- 1003 • **VQAv2:** General visual question answering task, evaluated using the official VQA Accuracy
1004 metric.
- 1005 • **COCO Captions:** Open-ended image captioning task, evaluated using CIDEr-D metric.
- 1006 • **DocVQA Layout:** Document QA task focusing on spatial layout and structure questions,
1007 evaluated using ANLS \times 100.
- 1008 • **DocVQA Number:** Document QA task focusing on numerical information extraction,
1009 evaluated using ANLS \times 100.
- 1010 • **ChartQA:** Chart interpretation and reasoning QA task, evaluated using the Relaxed Accu-
1011 racy metric.
- 1012 • **VTabFact:** Table reasoning multiple choice task, evaluated using accuracy.

1013 Results demonstrate consistent effectiveness across diverse real-world applications, with MeanShift
1014 achieving improvements in 15 out of 18 model-task combinations. The smaller improvement
1015 magnitudes compared to capability-focused benchmarks reflect the multi-faceted nature of these
1016 tasks, which require comprehensive reasoning abilities beyond isolated spatial or counting skills.

1017
1018
1019
1020
1021
1022
1023
1024
1025

1026 D DATASET EVALUATION DETAILS
1027

1028 In this section, we explain in detail how we prompt and evaluate the model’s performance across
1029 datasets and provide representative examples for each dataset. Each prompt consists of four com-
1030 ponents: **model prefix**, **task prefix**, **taxonomy prefix**, and **question**. The **model prefix**
1031 is the specific instruction token sequence required by different model families to perform certain
1032 tasks. For PaliGemma2 models, we use “answer en” as the model prefix, indicating that the model
1033 should answer in English for visual question answering tasks. For COCO dataset specifically, we use
1034 “caption en”, indicating that it is a captioning task. For Idefics3-8B-Llama3, no model prefix
1035 is required, so this component remains empty. The **task prefix** provides task-specific instructions
1036 that constrain the format of the model’s response. In multiple-choice questions, we use a task
1037 prefix such as “Answer the multiple choice question by only responding with the
1038 letter of the correct answer.” for example. In CLEVR and Super-CLEVR counting ques-
1039 tions, we use “Answer the question by only responding the number.” The **taxonomy**
1040 **prefix** of each taxonomy is the prompt we sampled in Section A.2, and it is only non-empty
1041 for the Prompt method. The **question** component contains the original question format from the
1042 dataset. Below are examples illustrating our prompting approach for each dataset.
1043



1070 Figure 10: Example prompt for the CV-Bench Relation dataset.
1071
1072
1073
1074
1075
1076
1077
1078
1079

1080
1081
1082 CV-BENCH COUNT
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102

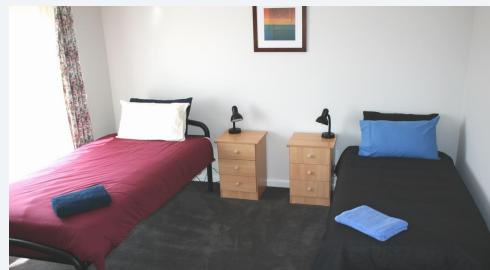


Image:

[Model Prefix] answer en [Task Prefix] Answer the multiple choice question by only responding the letter of the correct answer.
[Taxonomy Prefix] Prioritize counting objects and quantifying elements over other analysis. [Question] Answer the multiple choice question by only responding the letter of the correct answer. How many beds are in the image? Select from the following choices.

- (A) 0
- (B) 2
- (C) 1
- (D) 3
- (E) 4

1103
1104

Figure 11: Example prompt for the CV-Bench Count dataset.

1105
1106
1107
1108

1109 WHAT'SUP-A
1110



1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

Image:
[Model Prefix] answer en [Task Prefix] Answer the multiple choice question by only responding the letter of the correct answer.
[Taxonomy Prefix] Emphasize objects' positions relative to each other. [Question] Please select the correct caption for the image:
(A) A toilet roll under a chair
(B) A toilet roll to the left of a chair
(C) A toilet roll to the right of a chair
(D) A toilet roll on a chair

Figure 12: Example prompt for the What'sUp-A dataset.

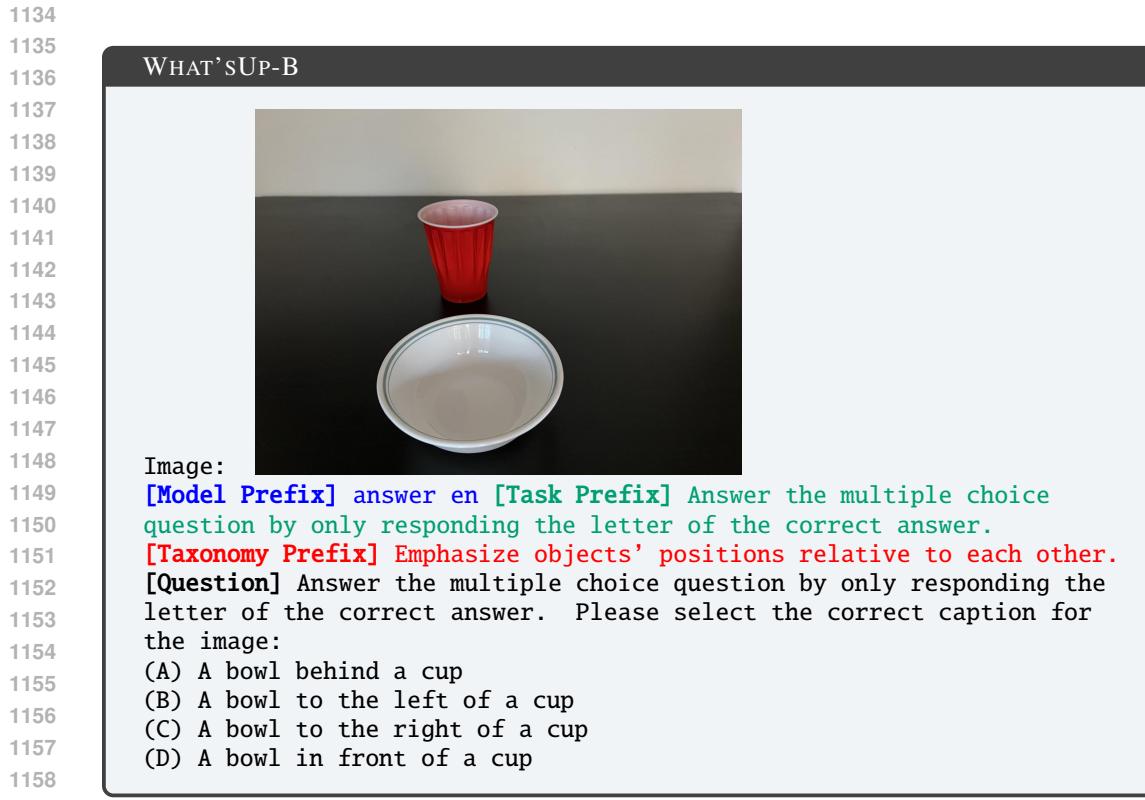


Figure 13: Example prompt for the What'sUp-B dataset.

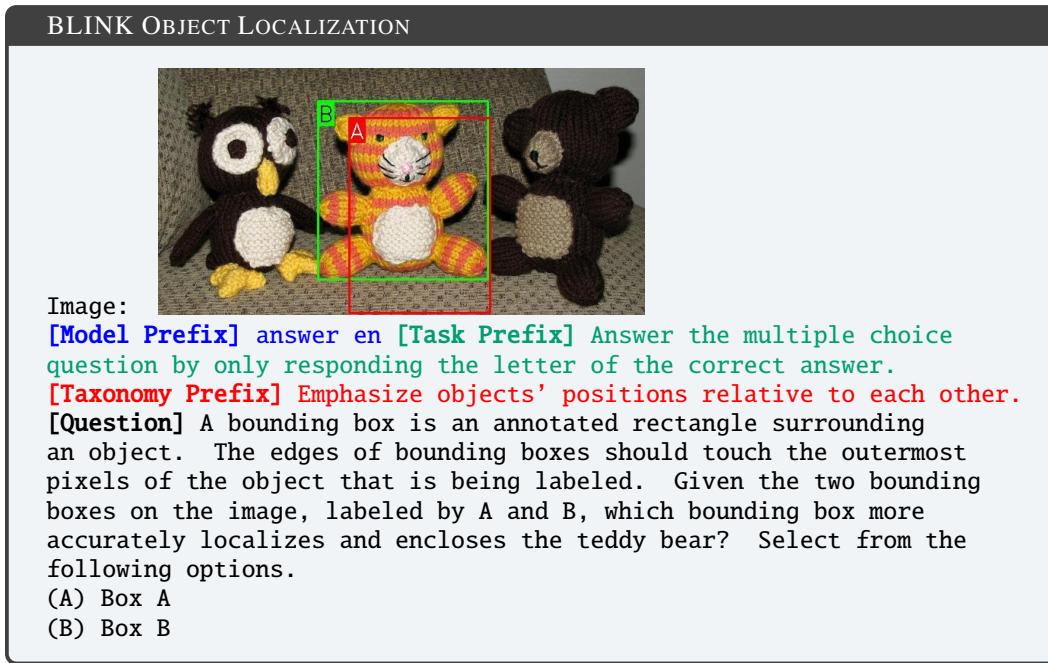
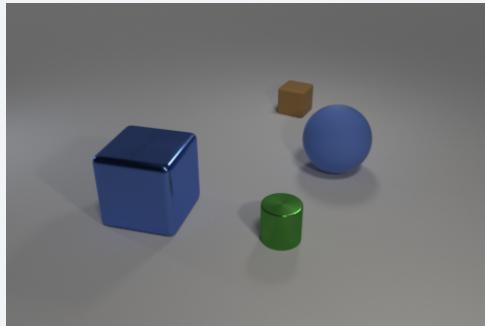
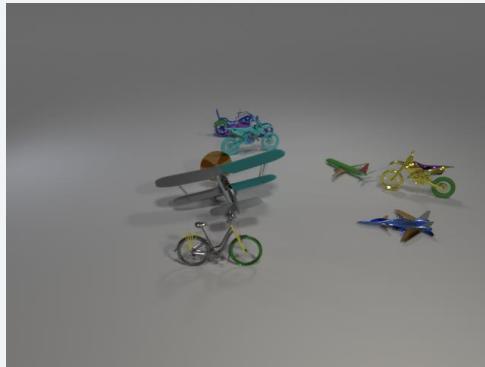


Figure 14: Example prompt for the BLINK Object Localization dataset.

1188
1189
1190
1191
1192 CLEVR
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203 Image:
1204 [Model Prefix] answer en [Task Prefix] Answer the question by only
1205 responding the number. [Taxonomy Prefix] Prioritize counting objects
1206 and quantifying elements over other analysis. [Question] How many
1207 different items are there in the image?
1208



1209
1210
1211
1212
1213
1214
1215
1216
1217
1218 SUPER-CLEVR
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231 Image:
1232 [Model Prefix] answer en [Task Prefix] Answer the question by only
1233 responding the number. [Taxonomy Prefix] Prioritize counting objects
1234 and quantifying elements over other analysis. [Question] How many
1235 different items are there in the image?
1236



1238 Figure 15: Example prompt for the CLEVR dataset.
1239
1240
1241



Figure 17: Example prompt for the VQAv2 dataset.



Figure 18: Example prompt for the COCO dataset.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307

1308 DocVQA LAYOUT
1309

1310
1311
1312
1313
1314
1315

1316
1317
1318
1319
1320
1321
1322
1323
1324

1325
1326
1327
1328
1329
1330

1331
1332
1333
1334
1335

1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349



Figure 19: Example prompt for the DocVQA Layout dataset.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361

1362 DocVQA Number
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384

1385 Image:
1386 [Model Prefix] answer en [Task Prefix] Answer the question about
1387 the image. Provide a short, direct answer. [Taxonomy Prefix]
1388 Prioritize counting objects and quantifying elements over other
1389 analysis. [Question] How many nomination committee meetings has S.
1390 Banerjee attended?
1391

1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

Report on Corporate Governance

Attendance at Nominations Committee Meetings during the financial year

Director	No. of meetings attended
Y.C. Deveshwar	2
A. Banerjee	2
S. Banerjee	2
A.V. Gopal Kumar	2
S.H. Noor	2
S.S. Mohur	1
D.K. Mihatra	NA
P.S. Ramaswami	2
S.S.N. Pillai [¶]	NA
M. Shukla [¶]	NA
K. Vaidyanathan	2

¶ Appointed Member w.e.f. 16th January, 2013.

V. SUSTAINABILITY COMMITTEE

The role of the Sustainability Committee is to review, monitor and provide strategic direction to the Company's sustainability practices towards fulfilling its triple bottom line objectives. The Committee seeks to guide the Company in integrating its social and environmental objectives with its business strategies.

Composition

The Sustainability Committee presently comprises the Chairman of the Company and six Non-Executive Directors, four of whom are Independent Directors. The Chairman of the Company is the Chairman of the Committee.

The names of the members of the Sustainability Committee, including its Chairman, are provided under the section 'Board of Directors and Committees' in the Report and Accounts.

[The structure, processes and practices of governance are designed to support effective management of multiple businesses while retaining focus on each one of them.]

Meetings and Attendance

During the financial year ended 31st March, 2013, three meetings of the Sustainability Committee were held, as follows:

Sl. No.	Date	Committee Strength	No. of Members present
1	30th April, 2012	6	6
2	2nd May, 2012	6	6
3	28th March, 2013	7	7

Attendance at Sustainability Committee Meetings during the financial year

Director	No. of meetings attended
Y.C. Deveshwar	3
S. Banerjee	3
H.G. Powell	3
A. Rajs	3
B. See	3
M. Shukla [¶]	1
B. Vaidyanathan	3

¶ Appointed Member w.e.f. 16th January, 2013.

CORPORATE MANAGEMENT COMMITTEE

The primary role of the Corporate Management Committee is strategic management of the Company's businesses within Board approved direction / framework.

Composition

The Corporate Management Committee comprises all the Executive Directors and six senior members of management. The Chairman of the Company is the Chairman of the Committee. The composition of the Corporate Management Committee is determined by the Board based on the recommendation of the Nominations Committee.

22 ITC Report and Accounts 2013
Source: <https://www.industrydocuments.ucsf.edu/docs/snbx0223>

Figure 20: Example prompt for the DocVQA Number dataset.

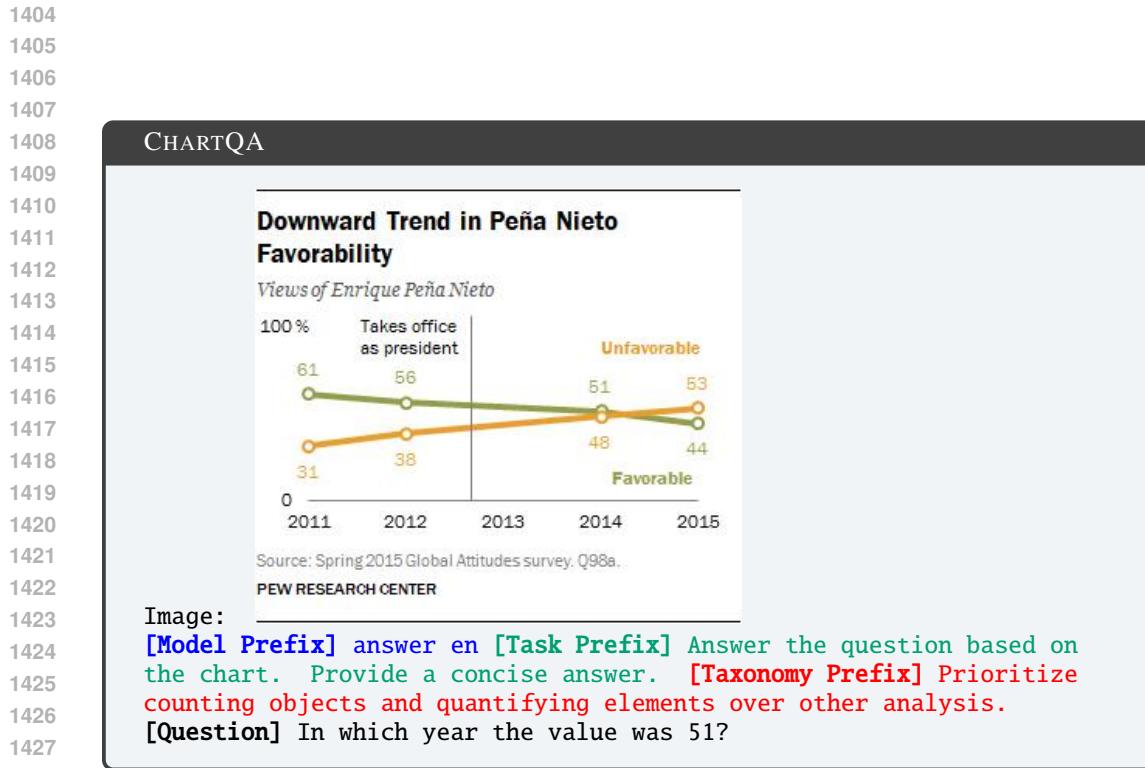


Figure 21: Example prompt for the ChartQA dataset.

1439
1440
1441
1442
1443
1444
1445 VTBFACT
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

team	head coach	years at school	overall record	record at school	acc record
boston college	frank spaziani	3	9 - 5	9 - 5	5 - 3
clemson	dabo swinney	3	13 - 8	13 - 8	9 - 4
duke	david cutcliffe	3	53 - 44	9 - 15	4 - 12
florida state	jimbo fisher	1	1 - 0	1 - 0	0 - 0
georgia tech	paul johnson	3	126 - 46	20 - 7	12 - 4
maryland	ralph friedgen	10	66 - 46	66 - 46	38 - 34
maine	randy shannon	4	21 - 17	21 - 17	11 - 13
north carolina	brian dawson	4	71 - 38	20 - 18	11 - 13
nc state	tom o'brien	4	60 - 66	16 - 21	9 - 15
virginia	mike London	1	24 - 5	0 - 0	0 - 0
virginia tech	frank beamer	24	229 - 115 - 4	187 - 92 - 2	38 - 10

Image:
[Model Prefix] answer en [Task Prefix] Answer the multiple choice question by only responding the letter of the correct answer.
[Taxonomy Prefix] Prioritize counting objects and quantifying elements over other analysis. [Question] ralph friedgen coach for 10 year at maryland
(A) Yes
(B) No

Figure 22: Example prompt for the VtabFact dataset.

1458 E LLM USAGE STATEMENT

1459

1460 LLMs were used for: (1) text polishing, (2) SAE feature verification using o3-mini (Section A.1), and
1461 (3) prompt generation using GPT-4o for baselines (Section A.2). These applications were limited to
1462 specific methodological components. Core research ideas, innovations, and conclusions are original
1463 author contributions.

1464

1465

1466

1467

1468

1469

1470

1471

1472

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511