

Anchored, Not Graded: Vision-Language Models Fail at Slant-from-Texture Perception [Preprint]

Qian Zhang¹, Michal Golovanevsky^{1,3}, Fulvio Domini², and James Tompkin¹

¹ Brown University Computer Science

² Brown University Cognitive Science

³ Harvard University

Abstract. Human perception of surface slant from texture exhibits systematic, graded biases that emerge reliably in psychophysical experiments. Prior work showed that unsupervised CNNs reproduce several human-like biases, while supervised CNNs do not. Do Vision-Language Models (VLMs) exhibit similar competences? Across multiple VLM families and model scales, zero-shot and in-context prompting both produce distinctive failures: slant is predicted at only a small set of anchors (e.g., 0° , $\pm 25^\circ$, $\pm 45^\circ$) with little dependence on stimulus field of view, optical slant, or surface curvature. Supervised fine-tuning partially remediates the failure, but residual anchoring persists. While success in high-level vision-language benchmarks might not require sensitivity to low-level geometric cues, we interpret anchoring as a failure at the representation-to-output language interface: Not necessarily an absence of geometric encoding, but a failure to express it in a graded form.

1 Introduction

Comparing human and neural network visual systems has a long anecdotal history, but recent large models trained on billions of images have made the comparison of their behaviors and mechanisms more meaningful. Understanding these relations through *in silico* experiments is an important direction for vision science, e.g., to raise hypotheses for how biological vision might function. But it is also useful for practical tasks. Often, human-AI collaboration is predicated on AI systems being able to predict estimates of human beings (AI: “*I can see that.*” Human: “*I cannot see that.*”). Yet, studies of whether the basic psychophysical abilities of AI systems match those of humans are rarer, as most works focus on downstream performance.

Take the simple task of estimating the slant of a surface. Humans are able to judge three-dimensional surface slant from texture gradients alone, yet these judgments are systematically biased rather than accurate. Psychophysical experiments using textured slant surfaces have identified at least three consistent biases in perception [34]: (B1) convex surfaces appear steeper than concave surfaces of equal physical slant; (B2) larger fields of view produce greater perceived slant; (B3) curvature-sign discrimination degrades at small fields of view, falling to chance at 5° FOV. These biases are commonly quantified by perceptual gain: the ratio of judged slant to ground-truth slant. This is approximately 0.56 for regular dot textures, indicating substantial underestimation of surface slant. Curvature-sign accuracy for such textures

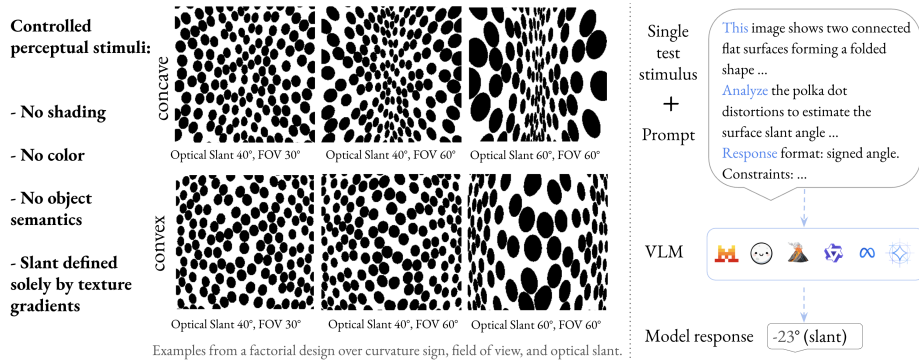


Fig. 1: Slant-from-texture as a controlled perceptual task for VLMs. *Left:* Synthetic dot-textured surfaces vary in curvature sign, field of view, and optical slant. Texture gradients are the only cue to 3D orientation. *Right:* Given a single image and instruction, we ask VLMs to estimate surface slant and curvature sign from texture alone.

averages around 86% overall but drops sharply at narrow fields of view. Thus, human slant-from-texture perception reflects systematic image-level heuristics, such as sensitivity to texture scaling gradients, rather than recovery of precise surface geometry.

Do neural networks exhibit similar biases? Wang et al. [40] investigated this question using convolutional neural networks (CNNs) trained on synthetic dot-textured stimuli. Unsupervised CNN autoencoders trained with a reconstruction objective reproduced human-like slant-from-texture biases. Human-like error patterns were recovered from the CNN’s internal representations using simple linear projections, which links representational structure to perceptual behavior. However, CNNs trained with direct supervision to regress physical slant achieved unbiased performance. This dissociation indicates that human-like biases arise not from architecture alone, but from learning objectives that emphasize texture statistics rather than explicit geometric labels. These results raise the question of whether such biases persist in visual representations both as models move beyond convolutional architectures and as they are trained under fundamentally different objectives.

Vision-Language Models (VLMs) present such a case. These models combine large-scale vision encoders, such as Vision Transformers, pretrained on broad image distributions, with language models trained on large-scale text corpora. Each is aligned through multimodal supervision on image-caption pairs. The architectural differences from CNNs and the distinct training routine is often assumed to yield representations that are richer or more semantically structured [7, 28, 46, 47]. Recent work has begun to question the geometric ability of VLMs in other domains [12, 32], but no studies have examined texture-based slant perception. Such stimuli contain no objects or semantic content; they contain only texture gradients that must be interpreted geometrically. Accordingly, despite strong performance on a range of vision-language benchmarks [2, 13, 17], it remains unknown whether VLMs exploit low-level texture cues in a manner comparable to human perception or unsupervised vision models and whether they can communicate slant in language.

We evaluate multiple VLM families on a classical slant-from-texture task, using stimuli and ground truth matched to prior human and CNN studies. We ask two

questions. First, behavioral: do VLMs exhibit the same systematic biases as humans and unsupervised CNNs, e.g, convex-concave asymmetry (B1) and field-of-view effects (B2, B3). We find that VLMs in zero-shot settings exhibit strong *anchoring* to a small set of discrete values, with predictions failing to exhibit monotonic dependence on stimulus parameters such as optical slant, field of view, or curvature sign—unlike the systematic biases observed in humans and unsupervised CNNs. Second, intervention: can standard adaptation techniques—prompt engineering, in-context learning with labeled examples, or supervised fine-tuning—induce such biases? We show that prompt engineering and in-context learning do not help. Supervised fine-tuning does introduce a monotonic relationship between predicted and ground-truth slant and improves curvature-sign discrimination, but does not eliminate anchoring: predictions remain clustered at discrete values rather than varying continuously. Probing the vision module within one VLM reveals a strong correlation between features and geometric quantities, suggesting that response anchoring is a language ‘readout’ problem. Together, these findings establish slant-from-texture as a diagnostic task for revealing systematic differences in how humans, unsupervised vision models, and language-supervised multimodal models respond to texture-based geometric cues, and show where effort is needed to produce future models that can predict (as needed) the human response to the world around us.

Assumptions and limitations. 1) Our experiments cover slant from random polka dot textures only, rather than a broader space of textures that includes regularity and shape variation [33] and perceptual judgements like depth or curvature magnitude. As psychophysical studies have shown a strong response in humans and CNNs to our polka dot slant stimuli, they are sufficient to establish VLM limitations. 2) We evaluate a breadth of models to characterize the landscape of VLMs, but specific alternatives may vary in significant ways. 3) Our findings do not obviate the usefulness of VLMs in end tasks or imply that VLMs must operate like humans; instead, they show differences in current VLM predictions with respect to human judgment that should be known when applying VLMs to tasks that require similar judgments.

2 Dataset and Experiment Setup

Stimuli. We produce synthetic dot-textured surfaces following Todd et al. [34, 35], and as followed by Wang et al. [40] (Fig. 1). These have no shading, color, or silhouette cues—slant must be perceived from texture gradients. Each stimuli depicts two textured flat surfaces forming a dihedral angle, varying along four parameters:

Optical slant σ_{cen}	10 values from 25° to 60°
Field of view (FOV)	10 values from 5° to 60°
Curvature sign	Convex or concave
Texture	12 i.i.d. random dot jitters (with different random seeds)

The surfaces have a true physical slant ρ that is their angle relative to the fronto-parallel plane. It is derived from σ_{cen} and FOV to preserve perceptual uniformity:

$$\rho = \sigma_{cen} - s \cdot \frac{\text{FOV}}{4}, \quad s = -1 \text{ (concave)}, \quad s = +1 \text{ (convex)}.$$

Optical slant σ_{cen} is the controlled experimental parameter matched across curvature conditions; physical slant ρ is what observers judge. Optical slant corresponds directly to the texture gradients while physical slant requires integrating curvature and FOV information. This makes physical slant prediction require geometric understanding of 3D structure rather than a direct ‘readout’ of local texture cues.

The first three attributes form a 200-condition factorial design (10 slant \times 10 FOV \times 2 curvature). As each instantiation of the 200 conditions has a random texture, each stimuli provides a new projection of the same geometry. Geometric reasoning requires an understanding of distortions to slant generalizable over appearance noise from texture jitter. We split the stimuli into two sets: 2000 for any training (fine-tuning, in-context learning, etc.) and 400 for testing.

Prompts. Each stimulus is paired with a text prompt. We vary the information provided in the prompt and its style to try to avoid any particular pitfall in how we describe the problem in language. This design allows us to test whether, and how, linguistic framing modulates visual estimations. The prompt varies the task, whether additional cues are given (e.g., “*The greater the slant angle, the more distorted the dots will appear.*”), and the style of prompt in its input and output. This produces seven key variants for physical slant angle prediction and three for binary concave/convex prediction (Tab. 1). We state all prompts in the supplementary material Tab. 4, Tab. 5, and Tab. 7.

Task type	Slant prediction (regression) vs. curvature sign (binary classification).
Instruction detail	Minimal vs. enriched with cue descriptions.
Language style	Natural (colloquial, e.g., ‘ridge/valley’, ‘fold’) vs. technical (e.g., ‘concave/convex’, ‘dihedral angle’).
Output format	Free text vs. structured JSON.
In-context	Absence vs. presence of example stimuli with slant labels.

Additionally, we provide an in-context learning prompt modifier. Along with the test stimuli, such prompts include four labeled example stimuli from the training set that are balanced across slant angle, field of view, and curvature. Labels are provided as free text (details in supplementary Tab. 6).

Table 1: Task prompt configurations. Natural language style prompt variants and their component. *Left:* Slant prediction as a regression task. *Right:* Curvature sign as a binary classification task.

Prompt Components Included		Prompt Components Included	
0	prompt_minimal_no_anchor	0	prompt_minimal
1	setup + task + format	1	setup + task + format
2	setup + task + format_eg	2	setup + task + cues + format
3	setup + task + format_eg_json		
4	setup + task + cues + format		
5	setup + task + cues + format_eg		
6	setup + task + cues + format_eg_json		

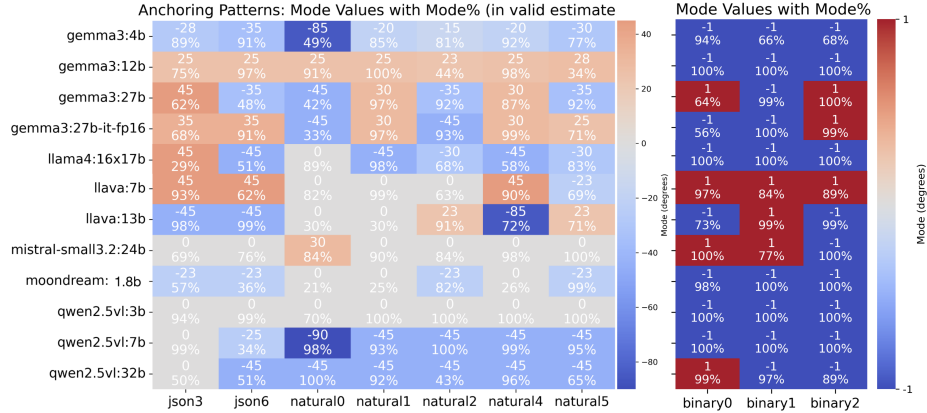


Fig. 2: VLMs anchor significantly on slant angle and sign prediction. Anchoring patterns for natural language prompts in regression task (left) and binary classification task (right). The heatmap shows distributions of estimated slant across VLMs (rows) and natural prompt variants (columns). All models collapse to a few discrete anchors (e.g., 0°, ±25°, ±45°) regardless of prompt detail, indicating that prompt variation does not mitigate anchoring. Prompt indices (e.g., ‘natural0’) and their components are in Tab. 4b and Tab. 7b, including both with and without JSON output constraints.

Model Families and Runs. We evaluated six recent open-sourced VLMs (Gemma3 [13], LLaMA4 [23], LLaVA [17], Mistral [18], Moondream [24], Qwen2.5-VL [27]) across multiple model parameter sizes, via local hosting using the Ollama API (v0.11.10) [25]. We use pretrained weights without fine-tuning unless specified. Qwen2.5-VL was evaluated using HuggingFace Transformers for SFT, probing, and attention-head ablation for weight access.

Each stimuli-prompt pair was submitted in prompt-completion format, and requests were executed asynchronously using multiprocessing. From the responses, we parsed numerical slant estimates or curvature sign classes. We balanced the number of concave and convex test stimuli across experiments. To capture trial-level variance, we repeated each query 10 times at the default temperature ($T = 0.7$). As the observed variance was negligible, subsequent experiments used a lower temperature ($T = 0.1$) with a single run per query for efficiency.

3 Results and Analysis

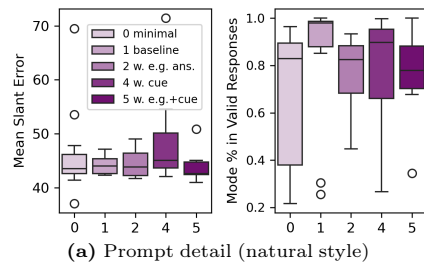
3.1 Zero-shot slant estimation

Systematic anchoring. Across all models evaluated, zero-shot VLM predictions exhibit large errors in slant and curvature sign (Fig. 2). Rather than varying smoothly with stimulus parameters, predicted slant values cluster around a small set of discrete anchor values, most commonly 0°, ±25°, and ±45°, regardless of prompt detail. The percentages below each cell indicate the mode (most frequent value) as a proportion of all responses. In many cases, the same value appears in more than 50% of responses, indicating anchoring rather than response noise.

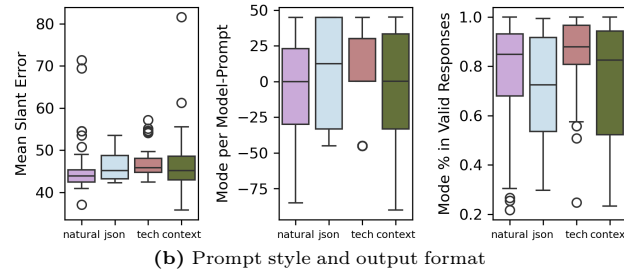
We observed three qualitative response patterns:

- *Instruction compliance errors*: Models ignored format constraints, outputting text without numerical values or producing out-of-range values.
- *Exemplar anchoring*: Prompts with examples induced “example anchoring” in some models, e.g., Moondream copying example prompt values verbatim (“-23°”), while other models were not affected by the specific values provided.
- *Zero-degree anchoring*: Models predicted flat surfaces (0°), sometimes accompanied by contradictory free-text reasoning in earlier runs (e.g., predicting flat while describing a steep slant).

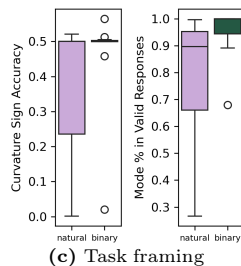
Anchoring persists across task framing, cues, prompt styles, and output formats. We compute the mean absolute error of slant estimates and variance for each prompt family and detail variants (Fig. 3). We observe no significant improvement in VLM ability to predict slant or sign. Varying model temperature did not improve performance. Trial-level consistency analyses at temperature 0.7 showed stable anchoring across repetitions with low variance; averaging predictions did not improve performance. Anchoring is a systematic response pattern rather than stochastic variability.



Varying prompt detail does not reduce anchoring. The anchor mode value changed for some models, but median errors are above 40° and anchoring dominates. Prompts emphasizing geometric cues yielded no consistent improvement (e.g., descriptions of texture scaling as included in prompts 4 and 5).



Prompts as natural vs. technical, or in JSON, or with context show no significant variation (two-way ANOVA on main effect of prompt family on slant error; Tab. 2). JSON formatting did not show consistent reduced example anchoring in any models but did exhibit increased overall variance, i.e. mode percentages decreased.



Curvature sign prediction as convex vs. concave (classification; ‘binary’ in plot) yielded no improvements over slant prediction (regression, ‘natural’ in plot), with accuracies around 0.5 (chance) and mode percentage increasing to around 90% in Fig. 3c. This trend is consistent across models.

Fig. 3: Effects of prompt and task framing as boxplots of 95% Confidence Intervals.

Source	sum_sq	df	F	p-value
C(model)	753.19	6	4.81	0.000291
C(prompt_type)	139.32	3	1.78	0.157403
Interaction	1230.32	18	2.62	0.001634
Residual	2192.87	84	-	-

Table 2: Two-way ANOVA (with interaction) of slant error mean for different models and prompts. Prompt type alone does not have a significant effect but depends on the model, e.g., one model might vary with JSON while another does not. Overall, there is no consistent main effect.

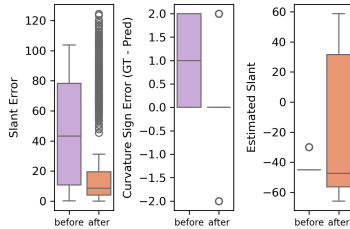


Fig. 4: Fine tuning only partially helps. Qwen2.5-VL before and after supervised fine-tuning: slant error decreases and anchoring weakens after SFT, but many outliers persist with high error.

In-context learning does not help. In-context prompts do not significantly differ from other prompt families (Fig. 3b): Median slant errors remain above 40° , curvature-sign accuracy stays near chance, and mode percentages remain high. ANOVA revealed no significant main effect of prompt type ($p=0.157$; Tab. 2). While specific anchor values sometimes shift across prompt conditions, the overall anchoring pattern persists.

Summary. In zero-shot settings, VLM predictions collapse to discrete anchors, show no systematic dependence on stimulus parameters, and are unaffected by prompt manipulations, including when given labeled examples for in-context learning. Model family and size modulate variability but do not alter the underlying anchoring pattern. Model performance is uniformly poor, making differences between models difficult to interpret. This behavior is consistent with evidence that LLM-style decoders exhibit anchoring effects and produce coarse numeric estimates that concentrate on a small set of preferred values (including round-number biases) [21, 37]. A complementary explanation is that round numbers are disproportionately frequent in natural language use, making a few canonical numerals strong attractors when models generate numbers as text [42].

3.2 Supervised fine-tuning (SFT)

We fine-tuned Qwen2.5-VL-3B using LoRA with a masked token loss on labeled slant and curvature data, updating the vision-language fusion layers (q_proj and v_proj).¹ Qwen2.5-VL-3B was chosen for its consistent anchoring patterns across regression and classification settings and its moderate size for efficient fine-tuning.

Improved slant estimation, some anchoring, but no low values. SFT substantially improves slant estimation and curvature-sign discrimination (Fig. 4). Mean slant

¹ LoRA was restricted to q/v_proj following standard PEFT practice; the regression-head result in Sec. 3.3 confirms this choice does not determine the bottleneck — frozen post-projector features are already linearly decodable for slant.

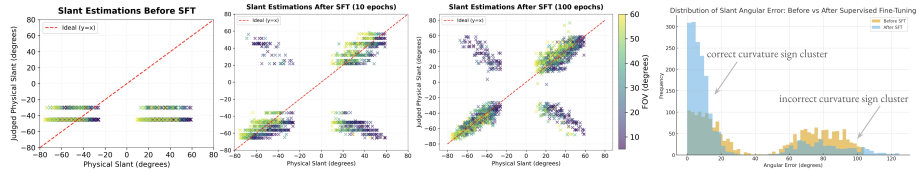


Fig. 5: SFT improves slant prediction but discrepancies remain. Judged vs. ground-truth slant angles. From left to right, top to bottom: VLMs before supervised fine-tuning show anchoring; bar chart of slant error distribution; VLMs with 10 epochs of SFT; VLMs with 100 epochs of SFT. After SFT, VLMs show improved slant prediction, but persistently estimate incorrect curvature sign and fail to predict low slants at all.

error decreases, and the distribution of predicted slant values broadens beyond the discrete anchors observed before fine-tuning. Curvature-sign accuracy improves from random chance (50%) to 86.10%. Despite these gains, anchoring is not eliminated. Compared to before fine-tuning, anchoring occurs at a finer scale and is distributed around the ground-truth slant (red line), forming horizontal bands near the red line of true prediction (Fig. 5). These decrease as training epochs increase, but remain present. Further, only high slant values are predicted $> \pm 25^\circ$.

Curvature sign remains in error. We use flip rates to analyze the proportion of trials in which the predicted curvature sign differs from ground truth (Sec. 3.4). Flip rates are elevated at small FOV and decrease systematically with increasing FOV and optical slant. These errors are *asymmetric* across curvature sign: concave stimuli show near-zero flip rates at large fields of view, whereas convex stimuli retain moderate flip rates (0.3–0.4) across optical slant. This pattern varies in the same direction as human behavior (B3), though is more pronounced. Paired statistical tests confirm that fine-tuning yields significant improvement ($df = 1999$, $p = 3.4 \times 10^{-106}$). Curvature-sign accuracy (86.10%) remains below that of unsupervised CNNs (96.4%) [40] and is numerically comparable to human performance (86%) [34], though the underlying error distributions differ qualitatively (Sec. 3.4).

Summary. SFT changes VLM behavior to partially remediate anchoring. Predictions become monotonically related to ground-truth slant (MAE: $45.1^\circ \rightarrow 15.3^\circ$; STD: $34.9^\circ \rightarrow 26.2^\circ$) and curvature-sign discrimination improves and becomes closer to past human and CNN studies (Fig. 10), though some slant errors and curvature-dependent asymmetries persist. This suggests that while SFT can steer model behavior, it does not fully override biases learned during pretraining.

Table 3: Layer-wise probe peaks across models.

Model	Vision tower	LM	Fusion	OS peak	FOV peak	PS peak	Curv peak
Qwen2.5-VL-3B	Qwen-native ViT	Qwen2.5	Late	0.998	0.921	0.879	0.995
LLaVA-1.5-7B	CLIP-ViT-L/14	Vicuna	Late (frozen vision)	0.997	0.923	0.905	0.990
PaliGemma-3B	SigLIP-So400M	Gemma-2B	Late (prefix)	0.997	0.921	0.926	1.000
Chameleon-7B	VQ-VAE	unified	Early (vocab fusion)	0.999	0.977	0.938	1.000

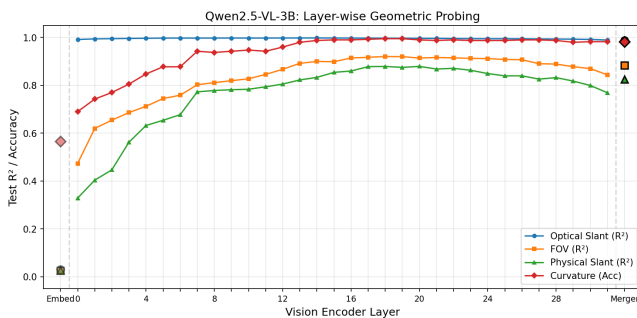


Fig. 6: Layer-wise linear probe performance across the 32-layer Qwen2.5-VL-3B vision encoder. Optical slant is near ceiling from the first transformer layer (actual $R^2 = 0.992\text{--}0.998$; appears flat at this scale). FOV, physical slant, and curvature require progressively deeper processing, peaking around layers 17–19 before declining in the final layers. Embed = Conv3D patch projection output (before any transformer processing). Merger = post-transformer spatial merging (2×2 tokens \rightarrow 2048-dim). Ridge regression ($\alpha = 1.0$) on mean-pooled 1280-dim features (2048-dim for merger).

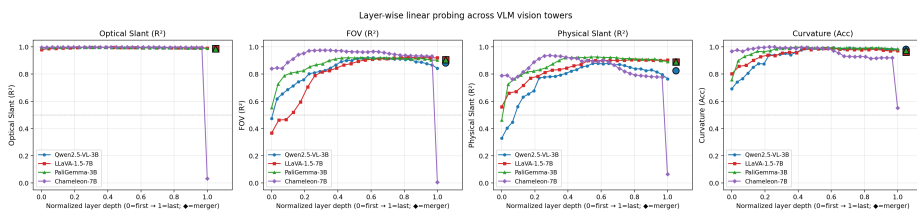


Fig. 7: Layerwise probing across four VLMs, extending the single-model analysis in Qwen2.5-VL-3B to LLaVA-1.5-7B (CLIP-ViT-L/14), PaliGemma-3B (SigLIP-So400M), and Chameleon-7B (VQ-VAE). The result is consistent with the encoding being architecture-general, with the readout bottleneck arising at the language interface rather than in the vision tower.

3.3 Probing the Vision Module

Pre-trained VLMs are trained on millions, often billions, of image–caption pairs, from which a complex structure is thought to emerge in their representations. Given that VLMs perform reasonably when directly trained on the stimuli and task, we use linear probing [1] to *localize* where correct behavior fails to emerge. The key question is whether zero-shot VLM failures arise from incorrectly encoded geometric information or from a deficit in language readout.

Probing setting. We analyze the vision encoder of Qwen2.5-VL-3B, which contains 32 transformer blocks (hidden size 1280, 10 heads) with RoPE and RMSNorm. Images are tokenized by a Conv3D patch embedding (14×14 patches) into 1280D tokens. A subsequent patch merger merges 2×2 tokens via a two-layer MLP, reducing token count by $4 \times$ and projecting to 2048D. For 224×224 stimuli, this produces 81 vision tokens (2048D) passed to the language model. To probe, we extract the 81 post-merger vision tokens in a zero-shot setting (no fine-tuning). Because feature dimensionality (2048) is comparable to sample count (2000), we

use ridge regression ($\alpha=1.0$) for all VLM probes to prevent overfitting; all reported R^2 values are evaluated on held-out test data.

Post-merger probing results. We evaluate regression fits using the coefficient of determination (R^2), which measures the proportion of variance in the target variable (PS, OS, FOV) explained by the probe from the latent tokens ($R^2=0$ indicates no explained variance; $R^2=1$ indicates perfect prediction). Ridge probes on mean-pooled VLM vision tokens achieve physical slant $R^2=0.826$, optical slant $R^2=0.988$, FOV $R^2=0.884$, and curvature accuracy $=0.983$, indicating that geometric information is strongly encoded in the vision representations. These results place the VLM vision encoder in the same performance regime as a pretrained ViT-MAE-86M (FOV $R^2=0.874$, PS 0.856, OS 0.997), suggesting VLM vision features capture geometry at a level comparable to this standalone vision backbone (a single comparison point).

Layer-wise probing. To trace how geometric information develops across the vision encoder, we extract mean-pooled features from every intermediate layer using forward hooks: the patch embedding output, all 32 transformer blocks, and the post-merger output. Applying ridge regression probes independently at each layer reveals a clear progression (Fig. 6). The raw patch embedding contains almost no geometric information (OS $R^2=0.031$, all others near zero; curvature accuracy $=0.565 \approx$ chance), indicating that the Conv3D projection captures little beyond local texture statistics. By the first transformer block (layer 0), however, optical slant is already near ceiling ($R^2=0.992$), confirming that it is a low-level image property recoverable from local texture gradients after minimal processing.

In contrast, FOV, physical slant, and curvature require progressively deeper processing: FOV R^2 rises from 0.474 (layer 0) to 0.920 (layer 18); physical slant from 0.329 to 0.879; curvature accuracy from 0.690 to 0.995. The correlation of all three variables peaks around layers 17–19 (the 56th–59th percentile of the 32-layer network) and then *declines* in the final layers: by layer 31, FOV drops to 0.844, physical slant to 0.769, and curvature to 0.983. The patch merger partially recovers from this late decline (FOV = 0.884, PS = 0.826), likely because its spatial 2×2 merging reintroduces local structure that the late layers had redistributed.

This late-layer decline is consistent with the view that the final transformer blocks specialize in the downstream language modeling objective, i.e., reorganizing features for cross-modal projection at the cost of some geometric accessibility. The practical implication is that middle layers of the vision encoder (around layer 18) contain the richest geometric signal, not the final output.

Generalization across vision encoders. To test whether the mid-/late-layer geometric-encoding peak is specific to Qwen2.5-VL’s vision tower, we replicate the analysis on three additional VLMs with different vision encoders and language models, Llava-1.5-7B, PaliGemma-3B, and Chameleon-7B², shown in Fig. 7. For each, we hook the per-layer outputs of the vision tower (including the projector), mean-pool

² As an architecture-diversity test, the sizes are each architecture’s canonical released size, not a controlled choice: LLaVA-1.5 has no 3B — smallest is 7B; PaliGemma ships in 3B only. A true size-controlled SFT comparison would be a separate experiment.

the patch tokens, and fit ridge regressions (FOV, optical slant, physical slant) and a logistic probe (curvature sign).

We found that geometric encoding generalizes across vision-encoder architectures. Across all four models, the layer-wise R^2 profile is qualitatively similar: optical slant saturates at $R^2 \geq 0.997$ by the first encoder layer; physical slant climbs gradually to a peak of 0.88-0.93 at 60-80% depth; curvature accuracy reaches 0.99-1.00; and the projector preserves $\geq 95\%$ of the late-layer peak. Our finding holds across vision towers.

One model of interest is Chameleon: the only model where vision and text tokens flow through the same transformer. Layer 31, the final transformer layer feeding into the LM head, shows a deep drop in geometric correlation: OS $R^2 = 0.033$, FOV $R^2 = 0.005$, PS $R^2 = 0.065$, curvature accuracy = 0.552 (chance). The geometric information present at layers 12-30 ($R^2 > 0.9$) is transformed away in the final layer toward the language-output distribution. This supports our claim that encoding is preserved deep in the network, but the final language-generation step discards it.

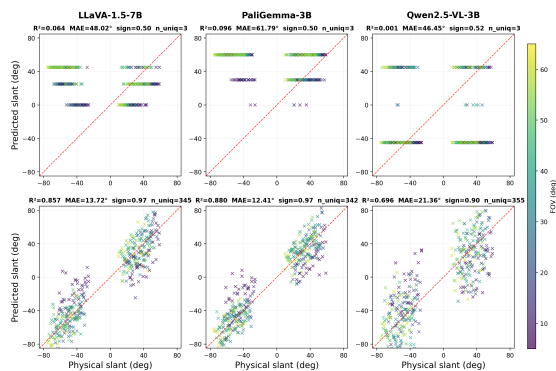


Fig. 8: Readout bottleneck - VLM output anchors (top); LM-input vision tokens encode continuous slant (bottom). We extract the mean-pooled post-projector vision tokens (in LM embedding space) from Qwen, LLaVa, and PaliGemma and train a linear regressor on train image tokens to predict physical slant. The figure shows results of the held-out test images.

Readout test. To test whether this is a true readout problem, we regress on the LM-input Vision Tokens, i.e. the vision tokens projected to LM embedding space contain the necessary geometric information but the language model cannot translate it into a continuous estimate. We mean-pooled the frozen post-projector tokens and trained a single ridge-regression head on the 2000-image training split, with no further tuning. As shown in Fig. 8, on the held-out 400-image test set, this single linear projection predicts physical slant at $R^2 = 0.696$ (Qwen2.5-VL-3B; MAE 21.4°), $R^2 = 0.857$ (LLaVA-1.5-7B; MAE 13.7°), $R^2 = 0.880$ (PaliGemma-3B; MAE 12.41°), and produces 342+ of 400 unique predicted angles, i.e., predictions

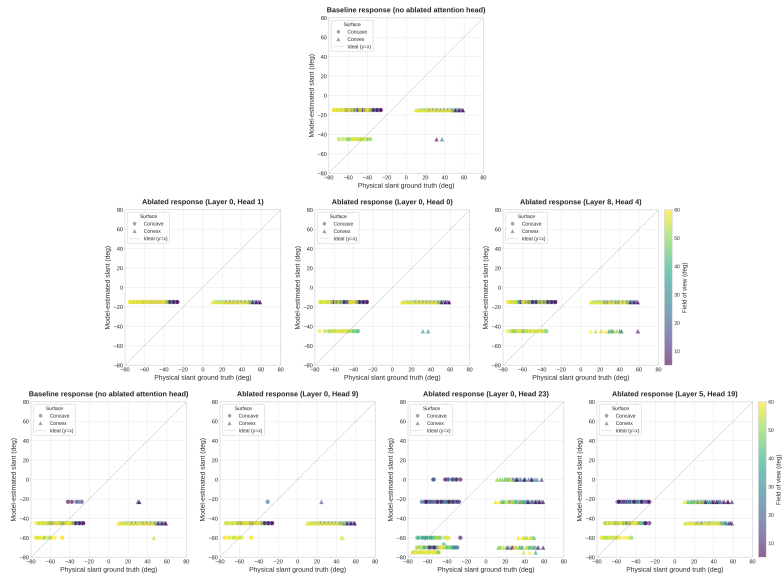


Fig. 9: Single attention head mean ablation on Qwen2.5-VL-3B (top) and Qwen2-VL-7B (bottom) slant prediction. (a) Baseline: predictions anchor at -45° . (b) Ablating L0 H9: predictions barely change. (c) Ablating L0 H23: anchoring spreads across more discrete values. (d) Ablating L5 H19: a new anchor at -23° appears that matches the prompt example (prompt copying). No ablation produces continuous stimulus-dependent predictions.

are essentially continuous. This shows directly the readout bottleneck, geometry is decodable, language head discards it.

3.4 Ablating Language Model Attention Heads

Recent work by Rudman et al. identified individual attention heads in VLMs that implement prompt-copying behavior, showing that ablating specific heads can eliminate copying of example values from the prompt [31]. Inspired by this approach, we test whether anchoring in our task arises from specific attention heads in the language model. For each targeted head, we replace its pre-projection activations with the mean activation across all batch and token positions, removing position- and content-specific information while preserving activation scale.

We sweep all layer-head pairs in Qwen2.5-VL-3B (36 layers \times 14 heads = 504 ablations), ablating one head at a time on the 400-stimuli test set. No single head ablation eliminates anchoring (Fig. 9). Some ablations modestly reduce the dominance of the primary anchor value and increase response variance, but most simply redistribute responses across existing anchor values without producing graded estimates that vary continuously with stimulus parameters.

We repeat the same analysis on the larger Qwen2-VL-7B model (28 layers \times 28 heads = 784 ablations) and observe a similar pattern: individual head ablations

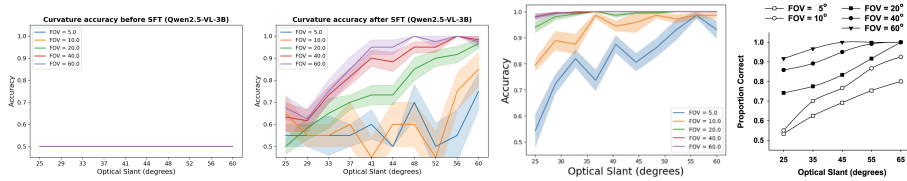


Fig. 10: VLM, CNN, and human curvature sign judgement accuracy versus optical slant, shown in matching axes for direct comparison. *Left:* Performance of VLM Qwen2.5-VL (3B) before supervised fine-tuning (SFT). *Center left:* After SFT. Qwen2.5-VL shows improved accuracy with increasing field of view, but still maintains high variance and is overall a poor proxy for human ability. *Center right:* Unsupervised CNN judgements reproduced from Wang et al. [40]. *Right:* Human judgements reproduced from Todd et al. [34].

slightly redistribute responses but do not substantially reduce anchoring. These results suggest that anchoring in continuous perceptual estimation is not attributable to a single attention head but instead reflects a distributed property of the language model. Although probing shows that geometric information is encoded in the vision representations, the language model cannot reliably read out continuous slant values. Consequently, attention-head interventions do not recover graded predictions.

4 Discussion

Anchoring as a failure of graded geometric expression. Our central finding is not that VLMs perform poorly on slant from texture, but that they fail in a specific, structured way. VLM predictions do not scatter randomly around the ground truth; they collapse to a small set of discrete values regardless of stimulus parameters. The anchoring pattern of VLMs is qualitatively distinct from the behavior of humans and unsupervised CNNs, which exhibit smooth, monotonic relationships between stimulus properties and perceived slant. Where humans show increasing perceived slant with larger fields of view (B2) and systematic differences between convex and concave surfaces (B1), VLMs show largely flat response curves—predictions that lack systematic variation with optical slant, FOV, or curvature sign. The failure is not one of magnitude or precision, but of graded expression: VLMs do not map continuous stimulus variation to continuous response variation.

Response anchoring is consistent with a language readout problem. We hypothesize that response anchoring in VLM slant judgments, i.e., the tendency to produce stereotyped numerical responses that do not co-vary with stimulus geometry, reflects a failure of the language model to read out this information. This is consistent with three observations:

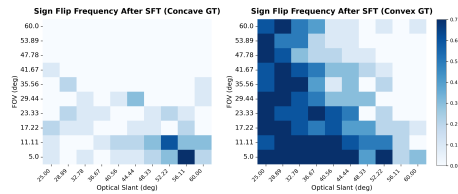


Fig. 11: SFT error analysis grid by slant angle and FOV, showing frequency of sign flips in predictions. Larger FOVs, larger optical slant images induced fewer flips; no curvature sign errors for concave FOV > 45°; convex-to-concave misjudge is more frequent than vice versa.

- SFT leaves the underlying inductive biases unchanged but improves how fused vision-language representations are translated into textual outputs, producing predictions that better follow the desired monotonic relation.
- Geometric information is present and linearly accessible throughout the vision encoder. Layer-wise linear probing shows that vision features strongly correlate with optical slant ($R^2 = 0.988$), curvature ($R^2 = 0.995$), and physical slant ($R^2 = 0.879$ at peak). Even the post-merger tokens passed to the language model retain high correlations with curvature ($R^2 = 0.983$) and physical slant ($R^2 = 0.826$).
- A standalone pretrained ViT (ViT-MAE-86M) with a comparable architecture achieves similar probe performance, suggesting vision encoding quality is not degraded—at least at intermediate layers—by VLM integration.

In summary, VLMs exhibit response-anchoring patterns that do not reflect the underlying geometric-encoding fidelity. In effect, *the model knows the geometry but does not say it*. Directly probing how geometric information transforms through language model hidden states remains future work.

Speculative: The cross-entropy bottleneck. A deeper question is whether the encoding-readout dissociation reflects a mismatch between training objectives. VLMs are trained with next-token prediction (cross-entropy over discrete text tokens): a continuous value like 42.5° must be expressed as a sequence of text tokens (“4”, “2”, “.”, “5”), and the model must learn that this token sequence corresponds to a specific point on a continuous scale. There is no direct encouragement for the model to produce text outputs that co-vary smoothly with the continuous geometric variables encoded in vision tokens.

Anchoring as a cue-deficient default value. Past human perception research found that distance prediction defaults to a specific scaling value when distance cues are not present [11]. Our stimuli are deprived of any cues beyond texture gradients. The models are trained on cue-rich images and defer to a small set of ‘default’ values. While this mechanism seems tenuous, given that geometric properties correlate to vision encoder features, future work could explore the impact of additional geometric cues within the stimuli on model prediction.

No low FOV predictions. Our stimuli are 224×224 . Image resolution, bit depth, and transformer tokenization block size likely affects the ability of the model to make predictions for low FOV stimuli because, at low FOV, texture gradients become small and so are ‘hard to see.’ This relates to ideas of human visual acuity.

Inference software variation. Specific anchoring values were not consistent across Ollama inference backends (e.g., slight differences in numeric execution), though the amount of anchoring did not vary notably.

Broader implications. Our findings suggest a tension between language grounding and continuous geometric expression in slant from texture task. VLMs perform well on high-level vision-language benchmarks that emphasize object recognition, scene understanding, and visual reasoning, yet they fail on a task that requires only the interpretation of low-level texture gradients. Standard VLM benchmarks emphasize categorical or semantic judgments that are naturally suited to discrete

text output, and may therefore not reveal difficulties in expressing continuous variables. Slant from texture, grounded in psychophysical research, provides a simple but revealing test case for probing the boundaries of elementary visual understanding in multimodal models, especially when considering tasks in which VLMs must accurately predict what a human is likely to see.

5 Related Work

Human perception of slant from texture. Slant-from-texture is a classic perceptual problem. Orccu et al. [26] compared linear perspective (a grid of lines) and texture gradient (diamond-shaped texture elements) cues for 75° slant perception. Observers tend to do better with combined cues than with either cue alone. The results were consistent with a linear combination of estimates from cues.

Todd et al. showed that human slant-from-texture perception is systematic and biased [34,35], identifying consistent effects of field of view and curvature sign (biases B1-B3 in our notation). More recent work further studied the effects of different texture cues in cue-conflicted and cue-consistent conditions. Chen et al. found that texture compression influenced slant settings more than other texture cues, with its influence decreasing with larger field of view (10° vs. 20°) and less regular textures [4].

Among texture types, regular blob (polka-dot) patterns are popular for modeling texture perception [16] and produce superior slant discrimination performance in humans compared to uniform lattices, Voronoi tessellations, plaids, and noise [30]. Unlike plaids or contour textures, they lack explicit perspective lines, making them a purer test of texture-gradient processing. They can also be systematically manipulated to test the effects of gradient compression [4] and gradient disruption [15].

Computational models of slant-from-texture perception. CNNs trained on texture statistics reproduce human-like slant biases under unsupervised learning objectives [40], supporting the use of deep networks as computational analogues for cue learning and inference (without claims to biological fidelity). This line of work connects to classic shape-from-texture formulations that recover surface orientation/shape from texture gradients under assumptions such as homogeneity [14,41]. Recent computer vision models revisit shape-from-texture under more realistic settings (unknown textures, nuisance factors) and show that texture alone can constrain local surface geometry in the wild [38]. In parallel, modern monocular depth estimation learns strong geometric priors from large-scale data, both supervised and self-supervised [8,10,43], but these systems typically entangle multiple cues and do not isolate the specific contribution of texture gradients to surface slant.

Vision-language models for visual perception. VLMs couple a pretrained vision encoder with a large language model via a lightweight alignment module, enabling strong performance on broad multimodal understanding and reasoning benchmarks [2,13,17,20,22,45]. However, their behavior on low-level perceptual and geometric tasks (e.g., fine-grained spatial relations, orientations, and viewpoint-dependent reasoning) is less well understood; existing spatially focused evaluations reveal substantial gaps [19,44]. Emerging diagnostics suggest that

large vision/VLM models exhibit weaknesses in metric depth and geometry understanding too, motivating targeted probes and RGB-D augmentation [3, 5, 6]. Overall, many training and evaluation approaches emphasize semantic recognition and language generation (captioning/VQA/grounding) rather than geometry.

Failure modes of VLMs. Controlled perceptual evaluations of VLMs remain rare. Zhang et al. [49] test five visual illusions (color constancy, color assimilation, color contrast, geometry relativity, and geometrical perspective), but most studies use naturalistic images. More broadly, a growing body of work documents systematic VLM failures, particularly on tasks requiring precise visual grounding and geometric reasoning. VLM failures on geometric primitives have been shown to depend upon spurious correlations rather than image evidence [9, 12, 29, 32, 36, 39, 48]. Mechanistic analyses suggest that some errors stem from over-reliance on language signals, including anchoring to prompt-provided cues and copying textual constraints even when they conflict with visual evidence [31]. However, much of this literature relies on naturalistic images, leaving open the question of how VLMs behave under tightly controlled psychophysical-style stimuli that isolate low-level geometric cues.

Zhang et al. [50] study image classification failures of VLMs. Their work provides tangential evidence for “underutilized visual features” in VLMs, considering classification where the output space is already discrete. We study a continuous geometric quantity, and find anchoring to a small set of salient values independent of stimulus parameters, which has no analog in discrete classification.

References

1. Alain, G., Bengio, Y.: Understanding intermediate layers using linear classifier probes. arXiv preprint arXiv:1610.01644 (2016) **9**
2. Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan, Y., Ge, W., Han, Y., Huang, F., et al.: Qwen technical report. arXiv preprint arXiv:2309.16609 (2023) **2, 15**
3. Cai, W., Ponomarenko, I., Yuan, J., Li, X., Yang, W., Dong, H., Zhao, B.: Spatialbot: Precise spatial understanding with vision language models. In: 2025 IEEE International Conference on Robotics and Automation (ICRA). pp. 9490–9498. IEEE (2025) **16**
4. Chen, Z., Saunders, J.A.: Multiple texture cues are integrated for perception of 3d slant from texture. *Journal of Vision* **20**(7), 14–14 (2020) **15**
5. Cheng, A.C., Yin, H., Fu, Y., Guo, Q., Yang, R., Kautz, J., Wang, X., Liu, S.: Spatialrgpt: Grounded spatial reasoning in vision-language models. *Advances in Neural Information Processing Systems* **37**, 135062–135093 (2024) **16**
6. Danier, D., Aygun, M., Li, C., Bilén, H., Mac Aodha, O.: Depthcues: Evaluating monocular depth perception in large vision models. In: *Proceedings of the Computer Vision and Pattern Recognition Conference*. pp. 20049–20059 (2025) **16**
7. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020) **2**
8. Eigen, D., Puhrsch, C., Fergus, R.: Depth map prediction from a single image using a multi-scale deep network. *Advances in neural information processing systems* **27** (2014) **15**

9. Gao, J., Pi, R., Zhang, J., Ye, J., Zhong, W., Wang, Y., Hong, L., Han, J., Xu, H., Li, Z., et al.: G-llava: Solving geometric problem with multi-modal large language model. arXiv preprint arXiv:2312.11370 (2023) 16
10. Godard, C., Mac Aodha, O., Firman, M., Brostow, G.J.: Digging into self-supervised monocular depth estimation. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 3828–3838 (2019) 15
11. Gogel, W.C.: The sensing of retinal size. *Vision research* **9**(9), 1079–1094 (1969) 14
12. Huang, K.H., Qin, C., Qiu, H., Laban, P., Joty, S., Xiong, C., Wu, C.S.: Why vision language models struggle with visual arithmetic? towards enhanced chart and geometry understanding. In: Findings of the Association for Computational Linguistics: ACL 2025. pp. 4830–4843 (2025) 2, 16
13. Kamath, A., Ferret, J., Pathak, S., Vieillard, N., Merhej, R., Perrin, S., Matejovicova, T., Ramé, A., Rivière, M., Rouillard, L., et al.: Gemma 3 technical report. arXiv preprint arXiv:2503.19786 **4** (2025) 2, 5, 15
14. Kanatani, K.i., Chou, T.C.: Shape from texture: General principle. *Artificial Intelligence* **38**(1), 1–48 (1989) 15
15. Kemp, J.T., Vishwanath, D., Domini, F.: Sensory uncertainty does not drive perceptual discriminability in 3d vision (2024) 15
16. Knill, D.C.: Surface orientation from texture: ideal observers, generic observers and the information content of texture cues. *Vision research* **38**(11), 1655–1682 (1998) 15
17. Li, F., Zhang, R., Zhang, H., Zhang, Y., Li, B., Li, W., Ma, Z., Li, C.: Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. arXiv preprint arXiv:2407.07895 (2024) 2, 5, 15
18. Liu, A.H., Khandelwal, K., Subramanian, S., Jouault, V., Rastogi, A., Sadé, A., Jeffares, A., Jiang, A., Cahill, A., Gavaudan, A., et al.: Ministral 3. arXiv preprint arXiv:2601.08584 (2026) 5
19. Liu, F., Emerson, G., Collier, N.: Visual spatial reasoning. *Transactions of the Association for Computational Linguistics* **11**, 635–651 (2023) 15
20. Liu, Y., Duan, H., Zhang, Y., Li, B., Zhang, S., Zhao, W., Yuan, Y., Wang, J., He, C., Liu, Z., et al.: Mmbench: Is your multi-modal model an all-around player? In: European conference on computer vision. pp. 216–233. Springer (2024) 15
21. Lou, J., Sun, Y.: Anchoring bias in large language models: An experimental study. *Journal of Computational Social Science* **9**(1), 11 (2026) 7
22. Lu, P., Bansal, H., Xia, T., Liu, J., Li, C., Hajishirzi, H., Cheng, H., Chang, K.W., Galley, M., Gao, J.: Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In: International Conference on Learning Representations (ICLR) (2024) 15
23. Meta: Llama 4 multimodal intelligence. Meta AI blog post, <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>, accessed: 2026-03-04 5
24. Moondream: Moondream. GitHub repository, <https://github.com/vikhyat/moondream>, accessed: 2026-03-04 5
25. Ollama: Ollama (software and documentation). Website, <https://docs.ollama.com/>, accessed: 2026-03-04 5
26. Oruç, I., Maloney, L.T., Landy, M.S.: Weighted linear cue combination with possibly correlated error. *Vision research* **43**(23), 2451–2468 (2003) 15
27. Qwen Team: Qwen2.5-vl technical report. arXiv (2025). <https://doi.org/10.48550/arXiv.2502.13923> 5
28. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PmLR (2021) 2

29. Rahmzadehgervi, P., Bolton, L., Taesiri, M.R., Nguyen, A.T.: Vision language models are blind. In: Proceedings of the Asian Conference on Computer Vision. pp. 18–34 (2024) [16](#)
30. Rosas, P., Wichmann, F.A., Wagemans, J.: Some observations on the effects of slant and texture type on slant-from-texture. *Vision Research* **44**(13), 1511–1535 (2004) [15](#)
31. Rudman, W., Golovanevsky, M., Arad, D., Belinkov, Y., Singh, R., Eickhoff, C., Mahowald, K.: Mechanisms of prompt-induced hallucination in vision-language models. arXiv preprint arXiv:2601.05201 (2026) [12](#), [16](#)
32. Rudman, W., Golovanevsky, M., Bar, A., Palit, V., LeCun, Y., Eickhoff, C., Singh, R.: Forgotten polygons: Multimodal large language models are shape-blind. In: Findings of the Association for Computational Linguistics: ACL 2025. pp. 11983–11998 (2025) [2](#), [16](#)
33. Todd, J.T., Thaler, L.: The perception of 3d shape from texture based on directional width gradients. *Journal of vision* **10**(5), 17–17 (2010) [3](#)
34. Todd, J.T., Thaler, L., Dijkstra, T.M.: The effects of field of view on the perception of 3d slant from texture. *Vision Research* **45**(12), 1501–1517 (2005) [1](#), [3](#), [8](#), [13](#), [15](#)
35. Todd, J.T., Thaler, L., Dijkstra, T.M., Koenderink, J.J., Kappers, A.M.: The effects of viewing angle, camera angle, and sign of surface curvature on the perception of three-dimensional shape from texture. *Journal of vision* **7**(12), 9–9 (2007) [3](#), [15](#)
36. Tong, S., Liu, Z., Zhai, Y., Ma, Y., LeCun, Y., Xie, S.: Eyes wide shut? exploring the visual shortcomings of multimodal llms. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 9568–9578 (2024) [16](#)
37. Tseng, C.Y., Roy, S., Thasin, M., Zhang, D., Effiong, B.: Streetmath: Study of llms’ approximation behaviors. arXiv preprint arXiv:2510.25776 (2025) [7](#)
38. Verbin, D., Zickler, T.: Toward a universal model for shape from texture. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 422–430 (2020) [15](#)
39. Vo, A., Nguyen, K.N., Taesiri, M.R., Dang, V.T., Nguyen, A.T., Kim, D.: Vision language models are biased. arXiv preprint arXiv:2505.23941 (2025) [16](#)
40. Wang, Y., Zhang, Q., Aubuchon, C., Kemp, J., Domini, F., Tompkin, J.: On human-like biases in convolutional neural networks for the perception of slant from texture. *ACM Transactions on Applied Perception* **20**(4), 1–18 (2023) [2](#), [3](#), [8](#), [13](#), [15](#), [22](#)
41. Witkin, A.P.: Recovering surface shape and orientation from texture. *Artificial intelligence* **17**(1-3), 17–45 (1981) [15](#)
42. Woodin, G., Winter, B., Littlemore, J., Perlman, M., Grieve, J.: Large-scale patterns of number use in spoken and written english. *Corpus Linguistics and Linguistic Theory* **20**(1), 123–152 (2024) [7](#)
43. Yang, L., Kang, B., Huang, Z., Xu, X., Feng, J., Zhao, H.: Depth anything: Unleashing the power of large-scale unlabeled data. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 10371–10381 (2024) [15](#)
44. Yu, S., Chen, Y., Ju, H., Jia, L., Zhang, F., Huang, S., Wu, Y., Cui, R., Ran, B., Zhang, Z., et al.: How far are vlms from visual spatial intelligence? a benchmark-driven perspective. arXiv preprint arXiv:2509.18905 (2025) [15](#)
45. Yue, X., Ni, Y., Zhang, K., Zheng, T., Liu, R., Zhang, G., Stevens, S., Jiang, D., Ren, W., Sun, Y., et al.: Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 9556–9567 (2024) [15](#)
46. Zhai, X., Wang, X., Mustafa, B., Steiner, A., Keysers, D., Kolesnikov, A., Beyer, L.: Lit: Zero-shot transfer with locked-image text tuning. 2022 ieee. In: CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 18102–18112 (2021) [2](#)

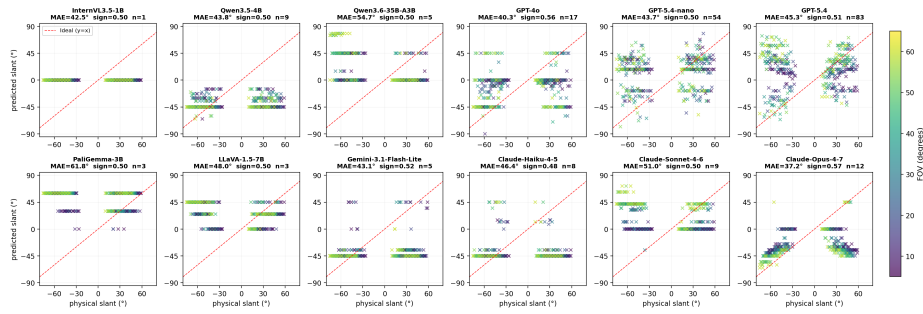


Fig. 12: Scatter plot of model predictions across 12 recent VLMs. Each panel shows predictions from a different model, with mean average error (MAE) for physical slant, accuracy for curvature sign discrimination, and the number of unique predicted values (out of 400 stimuli) in the title.

47. Zhang, J., Huang, J., Jin, S., Lu, S.: Vision-language models for vision tasks: A survey. *IEEE transactions on pattern analysis and machine intelligence* **46**(8), 5625–5644 (2024) [2](#)
48. Zhang, R., Jiang, D., Zhang, Y., Lin, H., Guo, Z., Qiu, P., Zhou, A., Lu, P., Chang, K.W., Qiao, Y., et al.: Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? In: *European Conference on Computer Vision*. pp. 169–186. Springer (2024) [16](#)
49. Zhang, Y., Pan, J., Zhou, Y., Pan, R., Chai, J.: Grounding visual illusions in language: Do vision-language models perceive illusions like humans? In: *Proceedings of Conference of Empirical Methods in Natural Language Processing*. EMNLP 2023 (2023) [16](#)
50. Zhang, Y., Unell, A., Wang, X., Ghosh, D., Su, Y., Schmidt, L., Yeung-Levy, S.: Why are visually-grounded language models bad at image classification? *Advances in Neural Information Processing Systems* **37**, 51727–51753 (2024) [16](#)

1 Prompt design and output formats

We include full prompt specifications to support reproducibility and to demonstrate that anchoring persists across substantial linguistic variation.

Full prompt text for each style of regression task: natural Tab. 4, technical Tab. 5, and in-context modifiers Tab. 6. SFT uses natural language prompts. Binary classification task also adopts natural language style prompts for consistency Tab. 7.

Note: earlier experiments included confidence ratings and free-text reasoning to assess response variability under stochastic decoding. As results showed that predicted slant values were strongly anchored and largely insensitive to decoding temperature, subsequent analyses focus on slant estimates alone.

2 Recent models

We also found anchoring in 12 frontier VLMs, including closed source ones. We evaluated GPT-4o, 5.4, 5.4-nano, Claude Opus 4.7, Sonnet 4.6, Haiku 4.5, Gemini

3.1 Flash-Lite, Qwen3.5-4B, Qwen3.6-35B-A3B (MoE), InternVL3.5-1B, LLaVA-1.5-7B, and PaliGemma-3B on the 400-stimulus test set with our natural-language prompt. *Every model* exhibits anchoring: their predictions collapse onto a few unique values vs. 200 distinct ground-truth physical-slant levels (Fig. 12).

For example, Qwen3.6-35B-A3B concentrates 274/400 predictions on exactly 0° ; Gemini-Flash-Lite places 301/400 at exactly -45° ; GPT-4o has 263/400 on just 0° and -45° . R^2 with physical slant is ≤ 0.16 for every model (mean 0.04), and sign-of-curvature accuracy is at chance (0.46–0.58). Thus, the anchoring bottleneck is not specific to open-weight or older VLMs, nor to a particular family (Google, Qwen, LLaVA, PaliGemma); all show the issue. GPT models shows a slightly different pattern with more variability and less extreme anchoring, suggesting some improvement. But slant estimation MAE remains high, and curvature sign discrimination is near chance.

3 Additional plots and anchoring analyses

In addition to natural language prompts (regression and binary classification), we include heatmaps for other prompt types.

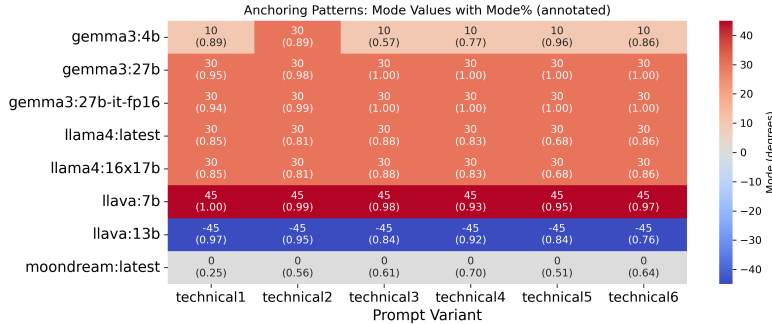


Fig. 13: Under technical prompt settings, mode value and frequency of results across models in heatmap.

4 Are failures due to out-of-distribution stimuli?

Our stimuli are synthetic polka-dot textures rendered under controlled conditions to isolate texture-gradient cues. A natural concern is that these images are out of distribution (OOD) relative to VLM training data, which consists predominantly of natural photographs and web content. If so, the observed anchoring failures might simply reflect OOD generalization failure — the model failing because the

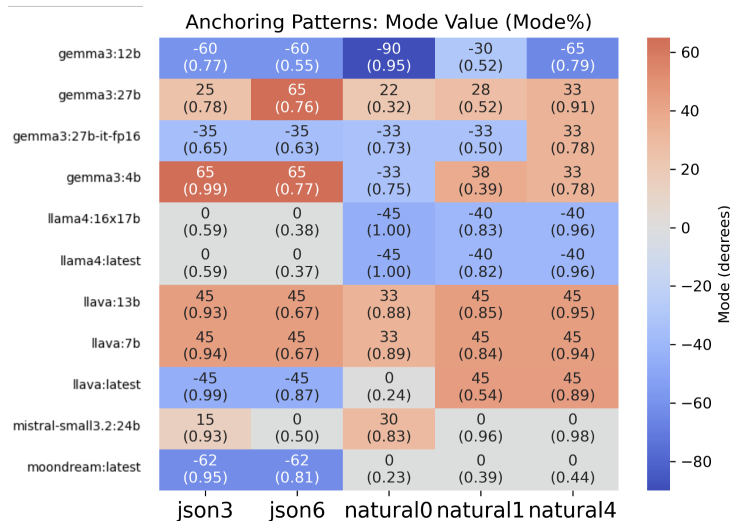


Fig. 14: With the in-context modifiers and natural language prompt settings, mode value and frequency of results across models in heatmap.

images are unlike anything in its training set — rather than a principled readout bottleneck at the language interface.

Several observations argue against a pure OOD account. First, polka-dot and regular blob textures appear pervasively in natural scenes and manufactured objects and are unlikely to be absent from large-scale web-crawled training corpora. Second, the failure pattern persists after supervised fine-tuning (SFT) on the task stimuli: even after training on thousands of polka-dot images, anchoring is not eliminated and sign-of-curvature accuracy remains below ceiling. An OOD explanation would predict that SFT, by familiarizing the model with the stimulus distribution, should largely remedy the failures. Third, sign-of-curvature accuracy is near chance even when curvature sign is the only required output, a binary discrimination that generic OOD degradation would not selectively impair. Finally, geometric probing confirms that the vision encoder extracts slant-relevant information from these stimuli at high fidelity (Sec. 3.3), ruling out a failure of low-level visual processing.

Taken together, these observations are consistent with the failures reflecting a readout bottleneck in the language model rather than an inability of the vision encoder to process the stimuli.

5 Inference backend variation: Ollama vs. HuggingFace

We ran a subset of experiments using two different inference backends for the same underlying model (Qwen2.5-VL-3B): Ollama and HuggingFace Transformers. The two backends produce numerically different anchor values—for example, Ollama may produce a dominant anchor at -45° while HuggingFace produces anchoring

at a different discrete value—likely due to differences in quantization, sampling defaults, or tokenizer handling between backends. Critically, however, the anchoring *pattern* is consistent: in both cases, predictions collapse to a small set of discrete values and do not co-vary systematically with stimulus parameters. All primary results reported in the main paper use HuggingFace Transformers for Qwen models, which provides direct access to model weights without additional quantization.

6 Attention head ablation details

Qwen2.5-VL-3B (36 layers \times 14 heads = 504 ablations). Baseline predictions collapse to two discrete anchor values with no prompt copying. Ablating individual heads merely redistributes which images anchor to which value without spreading the response distribution or changing the anchor values.

Qwen2-VL-7B (28 layers \times 28 heads = 784 ablations). Baseline predictions anchor at -45° with mode percentage 94% and correlation $r=0.149$ with ground-truth physical slant.

7 Metric notation: Wang et al. R vs. our r and R^2

Wang et al. use R to denote Pearson correlation; we use lowercase r for the same quantity. The two are equivalent; we adopt lowercase r to distinguish it from the coefficient of determination R^2 , which is a separate metric.

Wang et al. [40] report three distinct correlation metrics across their paper:

- **Wang et al. Table 1:** Pearson correlation R between the first two principal components of the latent space and physical variables (FOV, optical slant). These measure how well PCA axes align with individual variables.
- **Wang et al. Table 2:** Pseudo- R (not Pearson r) from a generalized linear model (GLM), with physical slant and FOV as independent variables and SVM latent distance as the dependent variable. Pseudo- R is a likelihood-ratio-based goodness-of-fit statistic, not equivalent to Pearson r or OLS R^2 . Reported values: 0.844 (concave) and 0.740 (convex).
- **Wang et al. Table 3:** Pearson correlation R between texture attributes (element length, width, area, spatial density) and the model’s “judged slant” (SVM latent distance).

For clarity of comparison, our R^2 is the coefficient of determination from ordinary least squares (OLS) regression: $R^2 = 1 - \text{SS}_{\text{res}}/\text{SS}_{\text{tot}}$. For the bivariate case, $R^2 = r^2$ (the square of Pearson r). Therefore, to compare our R^2 with Wang et al.’s Pearson R , one should square their values: e.g., their $R=0.819$ (Wang et al. Table 1, optical slant vs. 1st PC) corresponds to $R^2=0.671$.

The GLM pseudo R in Wang et al.’s Table 2 is not directly comparable to either Pearson r or OLS R^2 , as it is computed from a ratio of log-likelihoods rather than from residual variance. Therefore, we replicate their SVM-distance analysis on our models and report Pearson r for direct comparison.

Table 4: VLM prompts in the ‘natural language’ style, including variants and their components. Style changes in blue.

Components	Prompt
setup	The image shows two connected flat surfaces forming a folded shape, like a book or greeting card . Both surfaces are covered with identical round polka dots. The fold creates either a ‘ valley ’ (concave, opening away from you) or a ‘ ridge ’ (convex, pointing toward you).
task	Analyze the polka dot distortions to estimate the surface slant angle . IMPORTANT: Your angle estimate MUST be between -90 and +90 degrees . -90° = deepest possible valley (surfaces nearly horizontal pointing away), 0° = flat vertical surface, +90° = sharpest possible ridge (surfaces nearly horizontal pointing toward you). Also estimate your confidence between 0.0 and 1.0.
cues	Visual guide: Valley shapes compress and shrink dots near the fold center. Ridge shapes stretch and enlarge dots near the fold center. Stronger angles create more dramatic dot distortions.
format	Response format: (angle, confidence); brief_reasoning. Constraints: angle between -90 and +90, confidence between 0.0 and 1.0. Begin your response immediately with the opening parenthesis.
format_eg	Response format: (angle, confidence); brief_reasoning Example: (-23, 0.8); Dots are compressed near center indicating valley shape. Constraints: angle between -90 and +90, confidence between 0.0 and 1.0. Start your response immediately with the parentheses. No other text before it.
format_eg_json	Respond in this exact JSON format : {"angle": number, "confidence": number, "reasoning": "text"}. Constraints: angle between -90 and +90, confidence between 0.0 and 1.0. Example: {"angle": -23, "confidence": 0.8, "reasoning": "Dots compressed near fold"}
prompt_minimal no_anchor	Analyze this polka-dot folded surface. Valleys compress dots near the fold, ridges stretch them. Estimate slant angle (MUST be -90 to +90 degrees) and confidence (0.0-1.0). Response format: (angle, confidence); brief_reasoning. Begin immediately with the opening parenthesis.

(a) Prompt components and their details

Components	
0	prompt_minimal_no_anchor
1	setup + task + format
2	setup + task + format_eg
3	setup + task + format_eg_json
4	setup + task + cues + format
5	setup + task + cues + format_eg
6	setup + task + cues + format_eg_json

(b) Prompt variants and their component combinations.

Table 5: VLM prompts in the ‘technical language’ style, with variants and their components. Style changes in blue.

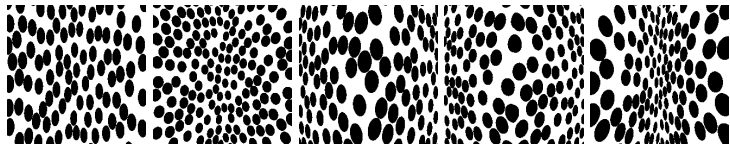
Components	Prompt
setup	The image shows a scene with two slanted, geometrically symmetric , flat surfaces textured with round polka dots of the same size. The slant is defined as the dihedral angle of the surfaces relative to the fronto-parallel plane . A convex shape points toward the viewer, while a concave shape points away.
task_signedslant	Please estimate the slant of the surface in degrees, within the range [-90, 90], where negative values indicate concave and positive values indicate convex.
format	IMPORTANT: Your answer MUST follow the format (slant_degree). Constraints: angle_number between -90 and +90. Begin your response immediately with the opening parenthesis. Do not include any extra text before or after your answer.
cues	Additional information: The slant will cause distortions in the polka dots, depending on both the sign and the magnitude of the slant angle. The greater the slant angle, the more distorted the dots will appear.
variations	The input image is from a set that includes both convex and concave shapes, with various surface slant angles and different fields of view.
effects	Near the intersection of the two planes, for a fixed field of view: the larger the convex angle, the larger and more stretched the dots appear. The larger the concave angle, the smaller and more compressed the dots appear. For a fixed slant angle, increasing the field of view increases the distortion of the dots.
biases	Previous psychophysical experiments have shown that humans consistently underestimate the slant of surfaces, especially for concave shapes. The accuracy of estimation is also affected by the field of view: smaller fields of view result in a larger underestimation bias.
hint	Recent research shows that when trained on the same stimulus images, unsupervised CNNs exhibit a similar bias to humans at test time, regardless of the neural network architecture. You may use this knowledge to adjust your estimation.

(a) Prompt components and their details.

	Components
1	setup + task_signedslant + format_signedslant
2	setup + task_signedslant + format_signedslant + cues
3	setup + task_signedslant + format_signedslant + cues + variations
4	setup + task_signedslant + format_signedslant + cues + variations + effects
5	setup + task_signedslant + format_signedslant + cues + variations + effects + biases
6	setup + task_signedslant + format_signedslant + cues + variations + effects + biases + hint

(b) Prompt variants and their component combinations.

Table 6: VLM prompt modifiers for in context settings. Style changes in blue.

In-context learning prompt	You are about to see a few example images and one test image (the last one). The true labels for the examples are as follows: (-62.778°, 1.0); (-49.444°, 1.0); (49.583°, 1.0); (37.222°, 1.0). Now here is the test image. This image shows ...				
	Example 1	Example 2	Example 3	Example 4	Test Image
					
GT labels	Concave, Physical slant -62.778°, FOV 11.11°	Concave, Physical slant -49.444°, FOV 35.56°	Convex, Physical slant 49.583°, FOV 41.67°	Convex, Physical slant 37.222°, FOV 60.00°	Concave, Physical slant 56.111°, FOV 53.89°

(a) Example in-context learning setup with labeled images and labels for few-shot slant estimation. The example images, test image, and prompts are sent to the model together in one query.

Table 7: VLM prompts for binary settings, with their variant components. Style changes in blue.

Components	Prompt
setup	Now here is the test image. This image shows two connected flat surfaces with polka dots, forming a folded shape. The fold creates either a ‘valley’ (like an open book) or a ‘ridge’ (like a roof peak).
task	Determine whether this shows a VALLEY or RIDGE.
visual_cues	Key visual indicators: In VALLEY configurations, polka dots near the fold appear smaller and more compressed. In RIDGE configurations, polka dots near the fold appear larger and more stretched.
format	Response format: (VALLEY/RIDGE). Begin your response immediately with the opening parenthesis.
cues	Visual guide: Valley shapes compress and shrink dots near the fold center. Ridge shapes stretch and enlarge dots near the fold center.
prompt_minimal	Look at this folded surface with polka dots. Is it folded like a VALLEY (inward fold, like an open book) or a RIDGE (outward fold, like a roof)? Response format: (VALLEY/RIDGE).

(a) Prompt components and their details

Components	
0	prompt_minimal
1	setup + task
2	setup + task + visual_cues

(b) Prompt variants and their component combinations