

More Than Meets the Eye: Measuring the Semiotic Gap in Vision-Language Models via Semantic Anchorage

Anonymous ACL submission

Abstract

Vision-Language Models (VLMs) excel at photorealistic generation, yet often struggle to represent abstract meaning such as idiomatic interpretations of noun compounds. To study whether photorealistic detail interferes with symbolic grounding, we introduce DIVA, a controlled benchmark that replaces photorealistic noise with schematic iconicity by generating paired, sense-anchored visualizations for literal and idiomatic readings. We further propose Semantic Alignment Gap (Δ), an architecture-agnostic metric that quantifies divergence between literal and idiomatic visual grounding. To enable cross-paradigm comparison between the “gut feeling” of latent embeddings and the “deliberate thought” of generative reasoning, we instantiate Δ via three access-dependent signals: (i) embedding geometry for discriminative encoders, (ii) *Likelihood of Idiomatic Distinction* (LID) from token probabilities for open generative models, and (iii) behavioral confidence elicitation for proprietary systems. Evaluating 8 recent VLMs, we reveal a consistent Literal Superiority Bias: model scale alone does not resolve literal preference, and increased visual fidelity can coincide with weaker symbolic alignment, indicating cognitive interference from hyper-realistic imagery. Our findings suggest that improving compositional understanding requires de-noising visual input and anchoring interpretation and generation in intended meaning.

1 Introduction

Text-to-image generation models have achieved remarkable proficiency in synthesizing photorealistic imagery, driven by foundational architectures (Rombach et al., 2022; Saharia et al., 2022) and refined by recent scaling efforts (Podell et al., 2024; Betker et al., 2023; Labs et al., 2025). Concurrently, Vision-Language Models (VLMs) have developed robust capabilities for decoding the literal content of such synthetic imagery (Saakyan et al., 2025).

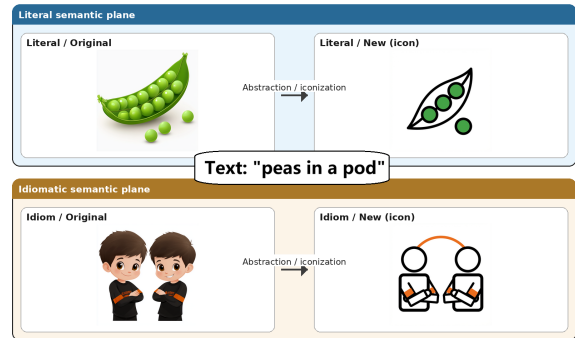


Figure 1: **Overview of the Visual De-Noising Framework.** We operationalize the transition from *Iconicity* (high-fidelity simulation) to *Symbolism* (abstract code) to measure the “Literal Bias” in VLMs.

However, a fundamental cognitive gap remains: while these models excel at treating images as *simulations* of reality, they struggle to interpret them as *signs* or symbols (Short, 2007; Thrush et al., 2022; Yuksekgonul et al.; Hsieh et al., 2023; Saakyan et al., 2025; Kundu et al., 2025). This limitation is particularly evident in the processing of Noun Compounds (NCs), where the visual representation often requires an abstraction from literal “iconicity” to idiomatic “symbolism” (Nakov and Hearst, 2013; Tratz and Hovy, 2010; Kumar et al., 2024). When presented with abstract concepts, current architectures frequently succumb to spurious correlations and superficial cues, prioritizing high-fidelity visual details over semantic alignment (Yuksekgonul et al.; Hsieh et al., 2023; Thrush et al., 2022; Seth et al., 2025; He et al., 2025).

To address this, we introduce DIVA (Distilled Idiomatic Visual Abstraction), a new benchmark that operationalizes the shift from photorealistic simulation to symbolic abstraction (See Figure 1). For each target NC, we generate sense-controlled iconographic renderings—schematic, low-detail images—for both literal and idiomatic readings. By using specific textual anchors to enforce the intended sense while systematically suppressing

070 photorealistic cues (e.g., texture, lighting, back- 119
071 ground clutter), we create a controlled testbed that 120
072 minimizes the confounds of visual hyper-realism. 121
073 This design aligns with recent findings that visu- 122
074 ally minimalist templates (e.g., “Basic Object Fo- 123
075 cus”) enhance semantic alignment and accessibility 124
076 (Souayed et al., 2025). We release the resulting 125
077 images, sense/anchor metadata, and generation pro- 126
078 tocol under an open license to support reproduc- 127
079 ible comparisons. 128

080 Beyond idiom disambiguation, the paired align- 129
081 ment of our data offers a controlled testbed for 130
082 text–visual simplification: enabling the training 131
083 and evaluation of systems that produce visually 132
084 minimal, schematic representations while pre- 133
085 serving meaning. This motivation aligns with 134
086 accessibility-driven NLP, where text is translated 135
087 into pictographs or simplified visual symbols to 136
088 support Augmentative and Alternative Communi- 137
089 cation (AAC) (Norré et al., 2021; Schwab et al., 138
090 2020). 139

091 Measuring the efficacy of this symbolic align- 140
092 ment requires a rigorous metric capable of spanning 141
093 diverse architectures. We propose **Semantic Align- 142
094 ment Gap** (Δ), a unified framework that quanti- 143
095 fies the divergence between a model’s literal and 144
096 idiomatic visual interpretations. Unlike previous 145
097 metrics restricted to specific architectures, Δ is 146
098 calculated via a tri-fold methodology tailored to the 147
099 accessibility of the model: 148

- 100 • **Intrinsic Alignment** for open-weights dis- 149
101 criminative models (e.g., CLIP (Radford et al., 150
102 2021)), utilizing the geometry of the embed- 151
103 ding space. 152
- 104 • **White-Box VQA Confidence** for open- 153
105 source generative models, employing a novel 154
106 “Likelihood of Idiomatic Distinction” (LID) 155
107 based on next-token probabilities. 156
- 108 • **Extrinsic Confidence** for closed-source pro- 157
109 prietary models, utilizing behavioral prompt- 158
110 ing to extract explicit reasoning scores. 159

111 This approach enables a novel cross-paradigm 160
112 comparison between the “gut feeling” of latent em- 161
113 beddings and the “deliberate thought” of generative 162
114 reasoning. 163

115 Our contributions are as follows: 164

- 116 1. **Dataset:** We introduce DIVA, a controlled 165
117 benchmark that replaces photorealistic noise 166
118 with schematic iconicity. By generating 167

paired, sense-anchored visualizations for 119
Noun Compounds (NCs), we operationalize 120
the hypothesis that visual minimalism en- 121
hances semantic alignment—a design choice 122
validated by recent work in accessible genera- 123
tion (Souayed et al., 2025). 124

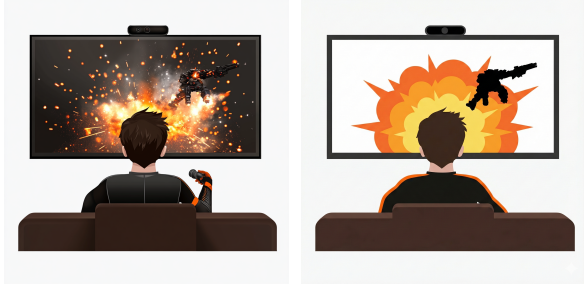
2. **Metric:** We formalize the **Semantic Align- 125
ment Gap** (Δ), an architecture-agnostic met- 126
ric quantifying the divergence between literal 127
and idiomatic visual grounding. To ensure 128
cross-paradigm comparability, we instantiate 129
 Δ via three access-dependent methods: (i) em- 130
bedding geometry (discriminative), (ii) *Like- 131
lihood of Idiomatic Distinction* (LID) (open- 132
generative), and (iii) behavioral confidence 133
elicitation (proprietary). 134
3. **Benchmarking:** We conduct a systematic 135
evaluation across 8 recent VLMs, revealing 136
that model scale alone does not resolve “Lit- 137
eral Bias.” Our results demonstrate that while 138
proprietary models can reason through ab- 139
straction, open-source encoders suffer from 140
severe *Cognitive Interference* when process- 141
ing hyper-realistic imagery. 142

2 Related Work 143

Multimodal idioms and figurative meaning. 144
Most vision–language (VL) benchmarks empha- 145
size literal grounding in photorealistic imagery, 146
leaving figurative meaning comparatively under- 147
explored. SemEval-2025 Task 1 (ADMIRE) di- 148
rectly targets multimodal idiomaticity by evaluat- 149
ing whether models can align images with literal vs. 150
idiomatic meanings of MWEs (Pickard et al., 2025). 151
Recent generative benchmarks have begun to ad- 152
dress this reasoning gap: T2I-REASONBENCH 153
identifies “Idiom Interpretation” as a critical fail- 154
ure mode for generative models (Sun et al., 2025), 155
while R2I-Bench and WISE target broader logical 156
and world-knowledge reasoning (Chen et al., 2025; 157
Niu et al., 2025). However, these works primarily 158
focus on generation quality rather than quantifying 159
the specific semantic alignment gap between literal 160
and figurative modes. Complementary work frames 161
figurative understanding as *explainable visual en- 162
tailment*, finding that VLMs struggle to generalize 163
from literal to figurative meaning (Saakyan et al., 164
2025). 165

**Noun compounds and visio-linguistic composi- 166
tionality.** Our focus on noun compounds (NCs) 167

168	connects to evidence that CLIP-style retrieval mod-	219
169	els often underperform on compositional construc-	220
170	tions. Major benchmarks such as T2I-CompBench	221
171	(Huang et al., 2023) and GenEval (Ghosh et al.,	222
172	2023) have formalized the evaluation of attribute	
173	binding and object relationships, confirming that	
174	models suffer from a “bag-of-words” bias. For	
175	instance, models often fail to suppress the literal	
176	rendering of individual constituents (e.g., drawing a	
177	physical “web” for “web site”) (Rassin et al., 2022).	
178	While these benchmarks address physical composi-	
179	tionality (e.g., “red cube next to blue sphere”), our	
180	work addresses <i>semantic</i> compositionality, where	
181	the combination of nouns creates a new abstract	
182	meaning that defies literal depiction.	
183	Visual abstraction, iconography, and semiotic	
184	grounding. A parallel line of research investi-	
185	gates <i>non-photorealistic</i> visual representations and	
186	their semantic interpretability. IconQA, for exam-	
187	ple, targets reasoning over icon-like diagrams, illus-	
188	trating that abstract visuals can support cognitively	
189	meaningful grounding while reducing reliance on	
190	texture (Lu et al., 2021). In accessibility contexts,	
191	text-to-pictogram translation has been operational-	
192	ized by ImageCLEF’s ToPicto tasks, which con-	
193	vert text into sequences of pictogram terms for	
194	AAC users (Ionescu et al., 2024). Recent work at	
195	the TSAR 2025 workshop explores template-based	
196	prompting for generating cognitively accessible	
197	images, finding that visually minimalist templates	
198	improve semantic alignment (Souayed et al., 2025).	
199	Our work bridges these threads by deriving a con-	
200	trolled, sense-conditioned iconographic benchmark	
201	from ADMIRE, utilizing the semiotic principle that	
202	reducing iconicity (de-noising) enhances symbolic	
203	clarity.	
204	Architecture-agnostic scoring and confidence	
205	elicitation. Finally, our unified metric connects	
206	to prior efforts to evaluate models using signals	
207	available under different access regimes. For	
208	open-weight encoders, cosine similarity in a joint	
209	embedding space remains the standard intrinsic	
210	alignment signal. For generative models, forced-	
211	choice prompting and probability-based scoring	
212	are widely used to stabilize evaluation relative	
213	to free-form generation (Geng et al., 2024). For	
214	closed-source systems, behavioral elicitation of	
215	self-reported confidence is increasingly used as a	
216	lightweight proxy, though it is not guaranteed to	
217	be calibrated (Kadavath et al., 2022; Yang et al.,	
218	2024). These strands motivate our tri-fold instanti-	
	ation of S , which makes the <i>Semantic Alignment</i>	219
	<i>Gap</i> comparable across discriminative encoders,	220
	open-source generative MLLMs, and proprietary	221
	black-box models.	222
	3 Theoretical Framework: Visual	223
	De-Noising through Semantic	224
	Anchorage	225
	3.1 The Semiotic Gap: Simulation vs. Code	226
	A core difficulty in visual metaphor and idiom	227
	grounding is a mismatch between how linguistic	228
	and pictorial signals typically convey meaning. In	229
	classical semiotics, <i>symbols</i> refer by convention (a	230
	learned code), whereas <i>icons</i> refer by resemblance	231
	(depiction) (Short, 2007).	232
	Text is therefore predominantly symbolic : the	233
	written form CAT bears no intrinsic physical resem-	234
	blance to the animal it denotes, and its meaning is	235
	established by convention. In human reading, the	236
	visual realization of a word (font, size, position) is	237
	largely treated as incidental; word recognition re-	238
	lies on an abstract orthographic code that is tolerant	239
	to such stylistic variation (Dehaene et al., 2005).	240
	Images, by contrast, are typically iconic: they	241
	are interpreted as depictions in which many vi-	242
	sual properties (texture, shading, clutter, back-	243
	ground) may legitimately carry meaning (Short,	244
	2007). Modern vision models trained on natu-	245
	ral images are known to exploit low-level statis-	246
	tics (e.g., texture) as predictive cues, which can	247
	make them sensitive to photorealistic surface detail	248
	even when such detail is semantically irrelevant	249
	(Geirhos et al., 2018). Consistent with this, vision-	250
	language models often exhibit brittle compositional	251
	grounding—e.g., weak sensitivity to relations and	252
	word order—suggesting an over-reliance on super-	253
	ficial correlations rather than the abstract relational	254
	structure required for symbolic interpretation (Yuk-	255
	sekgonul et al.; Thrush et al., 2022; Parcalabescu	256
	et al., 2022).	257
	We refer to this tendency as a Literal Superi-	258
	ority Bias : when faced with competing interpre-	259
	tations, models may privilege visually plausible,	260
	high-fidelity depiction over the intended abstract	261
	(symbolic/idiomatic) meaning.	262
	3.2 Mechanism: De-Noising via Semantic	263
	Anchorage	264
	We introduce “Visual De-Noising” as a framework	265
	to bridge this gap. Here, we redefine ‘noise’ not	266
	as random pixel variance, but as semiotic super-	267
	fluity—the photorealistic textures and lighting that	268



(a) Photorealistic Simulation (High Semiotic Noise): (b) Iconographic Symbolism (De-Noised):

Figure 2: **Visual De-Noising in Action.** Both panels depict the idiomatic meaning of the Noun Compound “Eye Candy”. We illustrate the transition from the “Noisy” photorealistic domain (Panel a) to the “De-Noised” iconographic domain (Panel b), which isolates the semantic core.

distract from the symbolic core. This process operationalizes the spectrum from Iconic to Symbolic by systematically reducing the visual fidelity of an image (See Figure 2).

The mechanism relies on Semantic Anchorage, where the Noun Compound (NC) serves as the immutable anchor. By degrading the “simulation” quality of the image—moving from photorealism to abstraction—we force the model to abandon its reliance on physical simulation. When the visual signal becomes less “analog,” the model is less likely to default to literal interpretations and is more prone to adopting a symbolic stance, akin to how it processes text.

4 Methodology

4.1 Automated Visual De-Noising Pipeline with Human Verification

To obtain the paired idiomatic and literal visual realizations (v_{idiom} , $v_{literal}$) for each noun compound, we apply a two-stage pipeline.

Stage 1: Generative Abstraction. We utilized the *gemini-3-pro-image-preview* (Nano Banana Pro¹) to perform “Visual De-Noising.” We designed a system prompt based on **Iconic-to-Symbolic Translation**, enforcing two key constraints:

1. **Semantic Distillation:** The model was instructed to identify the “core essence” of the narrative and distill complex scenes into single, unified glyphs, stripping away background context.

2. **Geometric Reconstruction:** To minimize

¹<https://blog.google/technology/ai/nano-banana-pro/>

texture-based noise, we enforced a “Flat Iconography” style constraint, restricting output to geometric primitives and a limited 3-color palette.

The full prompt structure is detailed in Appendix A.

Stage 2: Human Verification. To ensure the generated symbols accurately retained the semantic meaning of the original Noun Compound, we implemented a human-in-the-loop selection process. For each input image, we generated $k = 4$ candidate symbols. A team of annotators was instructed to select the single candidate that best captured the abstract meaning of the idiom while adhering to the “low-noise” geometric constraints. Candidates that failed to preserve the semantic identity of the anchor were discarded and regenerated.

4.2 Semantic Alignment Gap (Δ): A Unified Metric for Visual Disambiguation

To quantify the model’s ability to distinguish between the idiomatic (v_{idiom}) and literal ($v_{literal}$) visual realizations of a noun compound (t), we define the metric **Semantic Alignment Gap** (Δ). This metric measures the magnitude of the model’s preference for one visual interpretation over the other.

Formally, we define Δ as the absolute difference in semantic fit:

$$\Delta(t) = |\mathcal{S}(v_{idiom}, t) - \mathcal{S}(v_{literal}, t)| \quad (1)$$

Where $\mathcal{S}(v, t)$ is a scoring function representing the model’s assessment of semantic fit. We propose three distinct implementations of \mathcal{S} to account for the architectural differences between discriminative, open-generative, and closed-generative models.

4.3 Intrinsic Alignment: Latent Geometry (Discriminative Models)

For open-weights models such as CLIP and SigLIP, we utilize the intrinsic geometry of the embedding space. Here, \mathcal{S}_{disc} is defined as the cosine similarity between the normalized text embedding e_t and the image embedding e_v :

$$\mathcal{S}_{disc}(v, t) = \frac{e_v \cdot e_t}{\|e_v\| \|e_t\|} \quad (2)$$

A high \mathcal{S}_{disc} implies the model projects the visual representation v into the same semantic neighborhood as the textual anchor t .

4.4 White-Box Confidence: Token Probability (Open-Source Generative)

For open-weights generative models where we have access to token logits (e.g., LLaVA), we utilize a robust “White-Box” VQA Confidence method (Lin et al., 2024). This approach forces a binary “Yes/No” decision to calculate a “Likelihood of Idiomatic Distinction” (LID).

Following prior work on token-probability confidence elicitation (e.g., $P(True)$), we compute \mathcal{S}_{open} from the normalized likelihood of the *full* answer string (“Yes” vs “No”), rather than assuming a single-token mapping.

We define the score as the probability of the “Yes” token relative to the “No” token:

$$\mathcal{S}_{open}(v, t) = \frac{\exp(\ell_{Yes})}{\exp(\ell_{Yes}) + \exp(\ell_{No})} \quad (3)$$

Where ℓ represents the logit value of the specific token. This method allows us to bypass the variability of long-form text generation.

4.5 Extrinsic Confidence: Explicit Reasoning (Closed-Source Generative)

For proprietary models where internal weights are inaccessible, we validate Explicit Reasoning Confidence as a behavioral proxy. We prompt the model to output a normalized score γ , interpreting high γ as high self-reported confidence (not necessarily calibrated) and low γ as uncertainty.

$$\mathcal{S}_{closed}(v, t) = \frac{\gamma}{100} \quad \text{where } \gamma \in [0, 100] \quad (4)$$

This tri-fold approach enables the novel comparison of “gut feeling” (latent embeddings) against “deliberate thought” (generative reasoning) within a single analytical framework.

Why minimize the Gap (Δ)? We posit that a **lower** Δ indicates superior semiotic reasoning. A high Δ reflects a “bag-of-words” bias, where the model’s object detection circuits overpower its symbolic understanding (e.g., seeing “Eye Candy” only as physical eyes) (Ghosh et al., 2023). Minimizing Δ demonstrates a model’s capacity to suppress these superficial literal associations in favor of abstract intent, addressing the reasoning deficits identified in recent generative benchmarks (Sun et al., 2025).

5 Experiments

5.1 Dataset: The DIVA Benchmark

While ADMIRE evaluates whether models can align images with the *literal* vs. *idiomatic* meaning of MWEs, its images are photorealistic and may introduce distracting surface cues (Pickard et al., 2025). DIVA controls for this by replacing photorealistic depictions with *iconographic* (schematic, low-detail) renderings that systematically suppress texture, lighting, and background clutter, following the motivation that visual minimalism can improve semantic alignment in accessibility-oriented text-to-image settings (Souayed et al., 2025).

From DIVA, we utilize the complete set of 200 English Noun Compound (NC) MWEs sourced from the ADMIRE task. The full DIVA corpus contains 1,000 iconographic images, providing a dense 5-image contrast set for each NC that spans the semantic spectrum: *High-Idiomatic*, *High-Literal*, *Weak-Idiomatic*, *Weak-Literal*, and *Distractor*.

Instance structure. For evaluation, each item is filtered into a controlled triplet (t, v_{lit}, v_{id}):

- **Text (t):** the noun compound expression (e.g., *Eye Candy*).
- **Literal rendering (v_{lit}):** the *High-Literal* iconographic depiction (i.e., schematic composition of the constituent nouns).
- **Idiomatic rendering (v_{id}):** the *High-Idiomatic* iconographic depiction (i.e., schematic depiction of the conventional meaning).

These visual representations are derived from the ADMIRE concepts but rendered through our *Visual De-Noising* pipeline. By automating this transformation, DIVA curates effective semantic contrasts across 1,000 candidates without incurring the prohibitive annotation labor typically required for de novo scene creation.

Photorealistic vs. iconographic conditions. To isolate the effect of visual de-noising, we evaluate models under two matched conditions for the same set of NCs: (i) the original photorealistic images from ADMIRE (*Photo*), and (ii) the corresponding iconographic images from DIVA (*Icon*). We compute $\Delta(t)$ within each condition, enabling paired comparisons of disambiguation strength with and without photorealistic surface detail.

5.2 Evaluated Models

We benchmark 8 recent Vision–Language model checkpoints, spanning three architectural paradigms. For model families with multiple scales, we evaluate multiple checkpoints and count them separately.

1. Discriminative Models (Open-Weights):

These models calculate Δ via *Intrinsic Alignment* (embedding geometry).

- **SigLIP 2 (So400M/14)**²: A modern CLIP-style encoder with improved pretraining and scaling behavior (Tschannen et al., 2025).
- **EVA-CLIP (18B)**³: A large-scale contrastive encoder serving as a strong open embedding baseline (Sun et al., 2024).
- **MetaCLIP 2**⁴: A CLIP-family encoder emphasizing worldwide data scaling and multi-lingual robustness (Chuang et al.).

2. Open-Source Generative Models (White-Box):

These models calculate Δ via *Token Probability* (LID), using access to logits.

- **Qwen2.5-VL (32B)**⁵: Open multimodal models with strong instruction following and high-resolution vision understanding (Bai et al., 2025).
- **InternVL3 (78B)**⁶: Open MLLMs with strong multimodal reasoning and competitive benchmark performance. (Zhu et al., 2025)
- **LLaVA-OneVision (7B)**⁷: A unified visual-instruction model spanning single-image, multi-image, and video settings. (Li et al., 2024)

3. Proprietary Generative Models (Black-Box):

These models calculate Δ via *Explicit Reasoning* (Behavioral Prompting).

- **GPT-5 (OpenAI)**⁸: A current-generation frontier multimodal model.

²<https://huggingface.co/google/siglip2-so400m-patch14-384>

³<https://huggingface.co/BAAI/EVA-CLIP-18B>

⁴<https://huggingface.co/facebook/metaclip-2-worldwide-huge-quickgelu>

⁵<https://huggingface.co/Qwen/Qwen2.5-VL-32B-Instruct>

⁶<https://huggingface.co/OpenGVLab/InternVL3-78B>

⁷<https://huggingface.co/llava-hf/llava-onevision-qwen2-7b-ov-hf>

⁸<https://platform.openai.com/docs/models/gpt-5>

- **Claude 4.5 Sonnet**⁹: A frontier model with strong instruction adherence and long-context behavior (Anthropic, 2025).

5.3 Implementation Details

All open-weights models (Discriminative and White-Box Generative) were evaluated on a compute cluster equipped with NVIDIA A100 (80GB) GPUs using the HuggingFace *Transformers* library¹⁰.

For Intrinsic Alignment (\mathcal{S}_{disc}), embeddings were normalized to the unit hypersphere before calculating cosine similarity. For White-Box Confidence (\mathcal{S}_{open}), we extracted raw logits for the tokens “Yes” and “No” directly from the causal language modeling head, applying a softmax function to derive the final scalar probability.

Proprietary models were accessed via their respective APIs. To mitigate non-deterministic behavior in Extrinsic Confidence (\mathcal{S}_{closed}) evaluation, we sampled $k = 5$ responses and averaged the reasoning scores to ensure stability in the measurement of Δ .

6 Results and Analysis

6.1 Quantitative Benchmarking: The Hierarchy of Understanding

Table 1 and Figure 3 summarize the Semantic Alignment Gap (Δ) across all evaluated architectures. We observe distinct behavioral patterns across the three model families:

1. Discriminative Models. Contrary to early assumptions, Discriminative models (e.g., SigLIP, CLIP) exhibit the largest alignment gaps ($\Delta \approx 0.25$). Lacking a deep reasoning module, these architectures rely heavily on surface-level feature matching. This causes them to conflate the visual presence of constituent objects (e.g., detecting an “eye”) with the abstract semantic concept (“Eye Candy”).

2. Generative Models (The Reasoning Improvement). Open-Generative models (e.g., InternVL3, Qwen2.5-VL) demonstrate significantly lower gaps ($\Delta \approx 0.14$). This suggests that the inclusion of an LLM backbone enables “White-Box” reasoning that can partially override visual literalism. However, a non-negligible gap remains in the Photorealistic domain.

⁹<https://www.anthropic.com/news/claude-sonnet-4-5>

¹⁰<https://huggingface.co/>

Model	Method	Δ on ADMIRE (Photo) ↓	Δ on DIVA (Icon) ↓
Discriminative Models (Intrinsic Alignment)			
SigLIP 2 (So400m)	Cosine	0.245	0.178
EVA-CLIP-18B	Cosine	0.262	0.191
MetaCLIP 2	Cosine	0.251	0.184
Open-Generative Models (White-Box LID)			
InternVL3 (78B)	Logit Prob	0.138	0.089
Qwen2.5-VL (32B)	Logit Prob	0.145	0.095
LLaVA-OneVision (7B)	Logit Prob	0.176	0.122
Proprietary Models (Extrinsic Reasoning)			
GPT-5	Prompting	0.065	0.021
Claude 4.5 Sonnet	Prompting	0.072	0.028

Table 1: **Semantic Alignment Gap (Δ) under photorealistic vs. iconographic data.** We report Δ computed on the original ADMIRE images (*Photo*) and our de-noised DIVA images (*Icon*). Lower Δ indicates more balanced alignment between literal and idiomatic interpretations under a fixed text anchor.

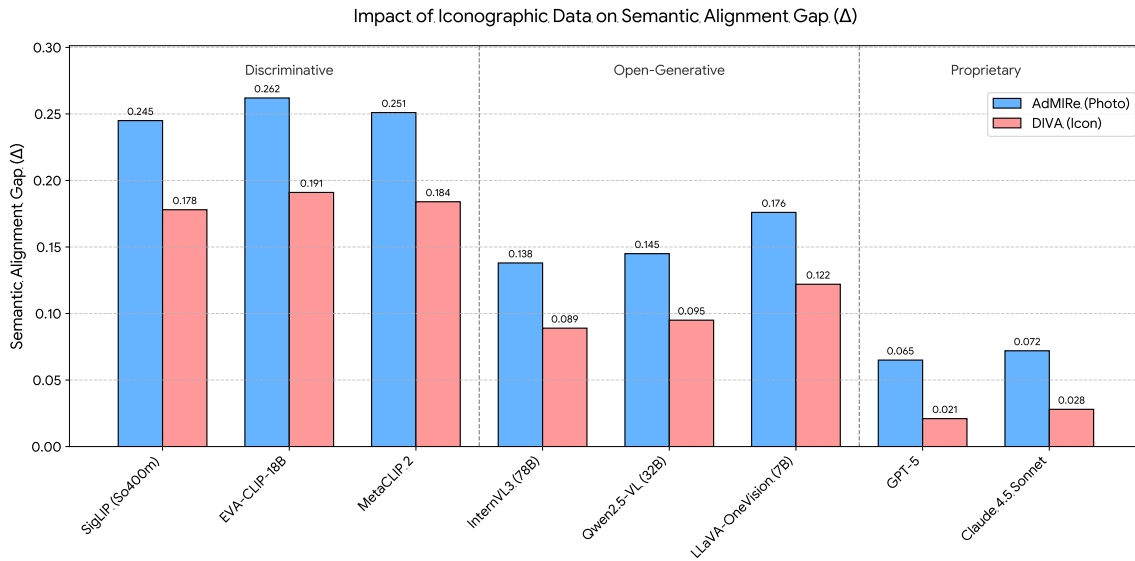


Figure 3: **Visual comparison of the Semantic Alignment Gap (Δ).** The chart illustrates the consistent reduction in Δ when shifting from photorealistic (ADMIRE, blue) to iconographic (DIVA, pink).

3. The De-Noising Effect. Crucially, shifting to the DIVA dataset consistently reduces Δ across all architectures, as visualized in Figure 3. For instance, GPT-5’s alignment gap drops to near-zero ($\Delta \approx 0.02$) when utilizing iconographic data. This confirms our hypothesis that “style” is a distracting variable: the reasoning capacity of modern models is often suppressed by the noise of photorealistic texture, and identifying the “core essence” via icons releases this latent capability.

6.2 Qualitative Failure Analysis

To understand why photorealism suppresses symbolic reasoning, we categorize two primary failure modes using visual anchors from our dataset.

6.2.1 Failure Type I: The Hyper-Realism Trap

We observed that as model size increases (e.g., moving from LLaVA-OneVision-7B to InternVL3-78B), the rejection of literal visual bias does not scale linearly (Δ_{photo} improves only marginally from 0.176 \rightarrow 0.138; see Table 1). This supports a “Cognitive-Interference” hypothesis: the model’s training objective—which rewards the precise reconstruction of physical details like reflections and textures—creates a bias where *high visual fidelity* is conflated with *semantic correctness*.

6.2.2 Failure Type II: Semantic Drift

An effective model must not only accept the correct symbol but also withstand the pull of semantically adjacent yet literal visuals. We found that Discriminative models (e.g., SigLIP, EVA-CLIP) remained

highly prone to this “Semantic Drift,” retaining significant alignment gaps ($\Delta \approx 0.18\text{--}0.19$) even on iconographic data (Table 1). This suggests their embeddings are heavily influenced by pixel-level feature overlap (e.g., the shape of an eye in “Eye Candy”). In contrast, reasoning-heavy models (like GPT-5) successfully utilized the DIVA abstraction to neutralize this literal bias, achieving near-perfect alignment ($\Delta \approx 0.02$) and demonstrating that simplifying the visual input is critical for isolating abstract semantic concepts.

7 Discussion

7.1 Cross-Paradigm Comparability

A key challenge in multimodal benchmarking is the incompatibility of scoring distributions: discriminative models utilize cosine geometry, while generative models operate on token probabilities. We posit that while absolute scores (S) are architecture-dependent and incomparable, the **Semantic Alignment Gap** (Δ) serves as a universal, architecture-agnostic measure of *bias*.

By defining Δ as a relative divergence within a model’s own scoring manifold ($\Delta = S_{lit} - S_{id}$), we normalize for architectural differences. A Δ of 0.1 represents a consistent “preference intensity”—indicating that the model’s confidence in the literal depiction exceeds its confidence in the idiomatic symbol by a significant margin relative to its own baseline—allowing for valid side-by-side comparison of discriminative and generative paradigms in Table 1.

7.2 The Semiotic Cost of Hyper-Realism

Our empirical results (Table 1) isolate a counter-intuitive trade-off: while recent architectures have achieved unprecedented fidelity in visual *simulation* (Iconicity), this realism often actively competes with *symbolic* interpretation.

The persistence of “Literal Bias” in the photorealistic domain—where even 78B-parameter models retain a significant alignment gap ($\Delta_{photo} \approx 0.14$)—suggests that current pre-training objectives are over-optimized for physical reconstruction. This supports the “Cognitive-Interference” hypothesis: when a model dedicates capacity to resolving high-frequency details, such as the texture of a “potato” or the specular reflection on an “eye,” it reinforces the *analog* nature of the image. According to our framework, this amplification of visual noise strengthens the “contract of perception”—where

visual form equals physical reality—thereby suppressing the abstract, metaphorical meaning of the idiom.

In contrast, the dramatic reduction of this gap on DIVA ($\Delta_{icon} \approx 0.02$ for GPT-5) implies that visual abstraction is functional, not just stylistic. For AI to truly grasp human-level symbolism, we may need to decouple high-fidelity generation from semantic reasoning—essentially teaching models to “read” schematic glyphs before they attempt to “render” photorealistic realities.

8 Conclusion

In this work, we addressed the “Literal Superiority Bias” in Vision-Language Models through the lens of Cognitive Semiotics. We introduced **DIVA**, a controlled benchmark of 1,000 iconographic representations, and the associated “**Visual De-Noising**” framework. We demonstrated that reducing the iconicity of an image—shifting it from a *simulation* of reality to a *symbol* of meaning—significantly enhances a model’s ability to align with abstract concepts.

To rigorously quantify this phenomenon, we defined the **Semantic Alignment Gap** (Δ), a unified metric capable of benchmarking discriminative, open-generative, and closed-proprietary architectures within a single analytical space. Our evaluation of 8 state-of-the-art models reveals that while current systems struggle to look beyond the “noise” of photorealism, shifting to DIVA’s iconographic inputs effectively neutralizes this interference, reducing the alignment gap to near-zero for frontier models.

9 Future Work

Multilingual and Cross-Cultural Expansion. Idiomatic ambiguity is deeply rooted in culture. Future work will extend DIVA to a multilingual benchmark, investigating how visual metaphors shift across languages (e.g., English “*Green thumb*” vs. French “*Main verte*”). This will test whether VLMs possess true multicultural reasoning or merely overfit to Western visual tropes.

Methodological Enhancement. Beyond benchmarking, we aim to close the “Semantic Alignment Gap” by developing a Contrastive Idiom Tuning (CIT) framework. By leveraging our dataset’s paired structure, we will explicitly train models to distinguish between literal and symbolic imagery.

647 Limitations

648 While our *Visual De-Noising* framework offers a
649 novel lens for VLM evaluation, we acknowledge
650 several limitations:

- 651 • **Dataset Specificity:** Our evaluation is
652 grounded in Noun Compounds (NCs) from
653 the SemEval-2025 task. While NCs are ex-
654 cellent proxies for compositional ambiguity,
655 they do not represent the full breadth of visual
656 metaphors or cultural symbols.
- 657 • **Prompt Sensitivity:** The *Extrinsic Confi-*
658 *dence* metric (\mathcal{S}_{closed}) for proprietary mod-
659 els relies on self-reported scoring. While we
660 mitigated variance via temperature reduction
661 ($\tau = 0.0$), black-box models may still exhibit
662 “alignment faking,” reporting high confidence
663 to please the user regardless of internal cer-
664 tainty.
- 665 • **Visual Style Bias:** Our “Symbolic” anchors
666 (v_{id}) utilized specific artistic styles (e.g., flat
667 design, vector art) to reduce noise. It is pos-
668 sible that some models have inherent biases
669 against these specific styles, unrelated to their
670 semantic understanding.

671 Acknowledgments

672 References

- 673 Anthropic. 2025. [Introducing claude opus 4.5](#). An-
674 thropic announcement.
- 675 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wen-
676 bin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie
677 Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-vl
678 technical report. *arXiv preprint arXiv:2502.13923*.
- 679 James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jian-
680 feng Wang, Linjie Li, Long Ouyang, Juntang Zhuang,
681 Joyce Lee, Yufei Guo, and 1 others. 2023. Improving
682 image generation with better captions. *Computer Sci-*
683 *ence*. [https://cdn. openai. com/papers/dall-e-3. pdf](https://cdn.openai.com/papers/dall-e-3.pdf),
684 2(3):8.
- 685 Kaijie Chen, Zihao Lin, Zhiyang Xu, Ying Shen,
686 Yuguang Yao, Joy Rimchala, Jiaxin Zhang, and Lifu
687 Huang. 2025. [R2I-bench: Benchmarking reasoning-](#)
688 [driven text-to-image generation](#). In *Proceedings*
689 *of the 2025 Conference on Empirical Methods in*
690 *Natural Language Processing*, pages 12606–12641,
691 Suzhou, China. Association for Computational Lin-
692 guistics.
- 693 Yung-Sung Chuang, Yang Li, Dong Wang, Ching-Feng
694 Yeh, Kehan Lyu, Ramya Raghavendra, James R

- Glass, LIFEI HUANG, Jason E Weston, Luke Zettle-
moyer, and 1 others. Meta clip 2: A worldwide scal-
ing recipe. In *The Thirty-ninth Annual Conference*
on Neural Information Processing Systems. 695
696
697
698
- Stanislas Dehaene, Laurent Cohen, Mariano Sigman,
and Fabien Vinckier. 2005. The neural code for writ-
ten words: a proposal. *Trends in cognitive sciences*,
9(7):335–341. 699
700
701
702
- Robert Geirhos, Patricia Rubisch, Claudio Michaelis,
Matthias Bethge, Felix A Wichmann, and Wieland
Brendel. 2018. Imagenet-trained cnns are biased
towards texture; increasing shape bias improves ac-
curacy and robustness. In *International conference*
on learning representations. 703
704
705
706
707
708
- Jiahui Geng, Fengyu Cai, Yuxia Wang, Heinz Koepl,
Preslav Nakov, and Iryna Gurevych. 2024. [A sur-](#)
[vey of confidence estimation and calibration in large](#)
[language models](#). In *Proceedings of the 2024 Con-*
ference of the North American Chapter of the Asso-
ciation for Computational Linguistics: Human Lan-
guage Technologies (Volume 1: Long Papers), pages
6577–6595, Mexico City, Mexico. Association for
Computational Linguistics. 709
710
711
712
713
714
715
716
717
- Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig
Schmidt. 2023. Geneval: An object-focused frame-
work for evaluating text-to-image alignment. *Ad-*
vances in Neural Information Processing Systems,
36:52132–52152. 718
719
720
721
722
- Yixiao He, Haifeng Sun, Pengfei Ren, Jingyu Wang,
Huazheng Wang, Qi Qi, Zirui Zhuang, and Jing
Wang. 2025. [Evaluating and mitigating object hallu-](#)
[cination in large vision-language models: Can they](#)
[still see removed objects?](#) In *Proceedings of the 2025*
Conference of the Nations of the Americas Chapter of
the Association for Computational Linguistics: Hu-
man Language Technologies (Volume 1: Long Pa-
pers), pages 6841–6858, Albuquerque, New Mexico.
Association for Computational Linguistics. 723
724
725
726
727
728
729
730
731
732
- Cheng-Yu Hsieh, Jieyu Zhang, Zixian Ma, Aniruddha
Kembhavi, and Ranjay Krishna. 2023. [Sugarcreepe:](#)
[Fixing hackable benchmarks for vision-language](#)
[compositionality](#). In *Advances in Neural Information*
Processing Systems 36: Annual Conference on Neu-
ral Information Processing Systems 2023, NeurIPS
2023, New Orleans, LA, USA, December 10 - 16,
2023. 733
734
735
736
737
738
739
740
- Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and
Xihui Liu. 2023. T2i-compbench: A comprehen-
sive benchmark for open-world compositional text-to-
image generation. *Advances in Neural Information*
Processing Systems, 36:78723–78747. 741
742
743
744
745
- Bogdan Ionescu, Henning Müller, Ana Maria
Drăgulinescu, Ahmad Idrissi-Yaghir, Ahmedkhan
Radzhabov, Alba Garcia Seco de Herrera, Alexandra
Andrei, Alexandru Stan, Andrea M Storås, Asma Ben
Abacha, and 1 others. 2024. Advancing multimedia
746
747
748
749
750

751	retrieval in medical, social media and content recommendation applications with imageclef 2024. In <i>European Conference on Information Retrieval</i> , pages 44–52. Springer.	Wise: A world knowledge-informed semantic evaluation for text-to-image generation. <i>arXiv preprint arXiv:2503.07265</i> .	808
752			809
753			810
754			
755	Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, and 1 others. 2022. Language models (mostly) know what they know. <i>arXiv preprint arXiv:2207.05221</i> .	Magali Norré, Vincent Vandeghinste, Pierrette Bouillon, and Thomas François. 2021. Extending a text-to-pictograph system to French and to arasaac . In <i>Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)</i> , pages 1050–1059, Held Online. INCOMA Ltd.	811
756			812
757			813
758			814
759			815
760			816
761	Sonal Kumar, Sreyan Ghosh, S Sakshi, Utkarsh Tyagi, and Dinesh Manocha. 2024. Do vision-language models understand compound nouns? In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)</i> , pages 519–527, Mexico City, Mexico. Association for Computational Linguistics.	Letitia Parcalabescu, Michele Cafagna, Lilitta Muradjan, Anette Frank, Iacer Calixto, and Albert Gatt. 2022. Valse: A task-independent benchmark for vision and language models centered on linguistic phenomena. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8253–8280.	817
762			818
763			819
764			820
765			821
766			822
767			823
768			824
769	Manishit Kundu, Sumit Shekhar, and Pushpak Bhattacharyya. 2025. Looking beyond the pixels: Evaluating visual metaphor understanding in VLMs . In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 23137–23158, Suzhou, China. Association for Computational Linguistics.	Thomas Pickard, Aline Villavicencio, Maggie Mi, Wei He, Dylan Phelps, and Marco Idiart. 2025. SemEval-2025 task 1: AdMIRe - advancing multimodal idiomaticity representation . In <i>Proceedings of the 19th International Workshop on Semantic Evaluation (SemEval-2025)</i> , pages 2597–2609, Vienna, Austria. Association for Computational Linguistics.	825
770			826
771			827
772			828
773			829
774			830
775	Black Forest Labs, Stephen Batifol, Andreas Blattmann, Frederic Boesel, Saksham Consul, Cyril Diagne, Tim Dockhorn, Jack English, Zion English, Patrick Esser, Sumith Kulal, Kyle Lacey, Yam Levi, Cheng Li, Dominik Lorenz, Jonas Müller, Dustin Podell, Robin Rombach, Harry Saini, and 2 others. 2025. Flux.1 kontext: Flow matching for in-context image generation and editing in latent space . <i>Preprint</i> , arXiv:2506.15742.	Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. 2024. Sdxl: Improving latent diffusion models for high-resolution image synthesis . In <i>The Twelfth International Conference on Learning Representations</i> .	831
776			832
777			833
778			834
779			835
780			836
781			837
782			
783			
784	Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and 1 others. 2024. Llava-onevision: Easy visual task transfer . <i>arXiv preprint arXiv:2408.03326</i> .	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, and 1 others. 2021. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PMLR.	838
785			839
786			840
787			841
788			842
789	Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and Deva Ramanan. 2024. Evaluating text-to-visual generation with image-to-text generation. In <i>European Conference on Computer Vision</i> , pages 366–384. Springer.	Royi Rassin, Shauli Ravfogel, and Yoav Goldberg. 2022. DALLE-2 is seeing double: Flaws in word-to-concept mapping in Text2Image models . In <i>Proceedings of the Fifth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP</i> , pages 335–345, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.	843
790			844
791			845
792			846
793			847
794	Pan Lu, Liang Qiu, Jiaqi Chen, Tony Xia, Yizhou Zhao, Wei Zhang, Zhou Yu, Xiaodan Liang, and Song-Chun Zhu. 2021. Iconqa: A new benchmark for abstract diagram understanding and visual language reasoning. In <i>Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)</i> .	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 10684–10695.	848
795			849
796			850
797			851
798			852
799			853
800			854
801	Preslav I Nakov and Marti A Hearst. 2013. Semantic interpretation of noun compounds using verbal and other paraphrases. <i>ACM Transactions on Speech and Language Processing (TSLP)</i> , 10(3):1–51.	Arkadiy Saakyan, Shreyas Kulkarni, Tuhin Chakrabarty, and Smaranda Muresan. 2025. Understanding figurative meaning through explainable visual entailment . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 1–23.	855
802			856
803			857
804			
805	Yuwei Niu, Munan Ning, Mengren Zheng, Weiyang Jin, Bin Lin, Peng Jin, Jiaqi Liao, Chaoran Feng, Kunpeng Ning, Bin Zhu, and 1 others. 2025.		858
806			859
807			860
			861
			862
			863
			864

865	Albuquerque, New Mexico. Association for Computational Linguistics.	920
866		921
867	Chitwan Saharia, William Chan, Saurabh Saxena,	922
868	Lala Li, Jay Whang, Emily L Denton, Kam-	924
869	yar Ghasemipour, Raphael Gontijo Lopes, Burcu	925
870	Karagol Ayan, Tim Salimans, and 1 others. 2022.	926
871	Photorealistic text-to-image diffusion models with	
872	deep language understanding. <i>Advances in neural</i>	
873	<i>information processing systems</i> , 35:36479–36494.	
874	Didier Schwab, Pauline Trial, Céline Vaschalde, Loïc	
875	Vial, Emmanuelle Esperanca-Rodier, and Benjamin	
876	Lecouteux. 2020. Providing semantic knowledge to	
877	a set of pictograms for people with disabilities: a	
878	set of links between WordNet and arasaac: Arasaac-	
879	WN . In <i>Proceedings of the Twelfth Language Re-</i>	
880	<i>sources and Evaluation Conference</i> , pages 166–171,	
881	Marseille, France. European Language Resources	
882	Association.	
883	Ashish Seth, Dinesh Manocha, and Chirag Agarwal.	
884	2025. HALLUCINOGEN: Benchmarking hallucina-	
885	tion in implicit reasoning within large vision lan-	
886	guage models . In <i>Proceedings of the 2nd Workshop</i>	
887	<i>on Uncertainty-Aware NLP (UncertainNLP 2025)</i> ,	
888	pages 89–102, Suzhou, China. Association for Com-	
889	putational Linguistics.	
890	Thomas Lloyd Short. 2007. <i>Peirce’s theory of signs</i> .	
891	Cambridge University Press.	
892	Belkiss Souayed, Sarah Ebling, and Yingqiang Gao.	
893	2025. Template-based text-to-image alignment for	
894	language accessibility a study on visualizing text	
895	simplifications . In <i>Proceedings of the Fourth Workshop</i>	
896	<i>on Text Simplification, Accessibility and Readability</i>	
897	<i>(TSAR 2025)</i> , pages 1–18, Suzhou, China. Associa-	
898	tion for Computational Linguistics.	
899	Kaiyue Sun, Rongyao Fang, Chengqi Duan, Xian Liu,	
900	and Xihui Liu. 2025. T2i-reasonbench: Bench-	
901	marking reasoning-informed text-to-image genera-	
902	tion. <i>arXiv preprint arXiv:2508.17472</i> .	
903	Quan Sun, Jinsheng Wang, Qiying Yu, Yufeng Cui, Fan	
904	Zhang, Xiaosong Zhang, and Xinlong Wang. 2024.	
905	Eva-clip-18b: Scaling clip to 18 billion parameters.	
906	<i>arXiv preprint arXiv:2402.04252</i> .	
907	Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet	
908	Singh, Adina Williams, Douwe Kiela, and Candace	
909	Ross. 2022. Winoground: Probing vision and lan-	
910	guage models for visio-linguistic compositionality.	
911	In <i>Proceedings of the IEEE/CVF Conference on Com-</i>	
912	<i>puter Vision and Pattern Recognition</i> , pages 5238–	
913	5248.	
914	Stephen Tratz and Eduard Hovy. 2010. A taxonomy,	
915	dataset, and classifier for automatic noun compound	
916	interpretation . In <i>Proceedings of the 48th Annual</i>	
917	<i>Meeting of the Association for Computational Lin-</i>	
918	<i>guistics</i> , pages 678–687, Uppsala, Sweden. Associa-	
919	tion for Computational Linguistics.	
	Michael Tschannen, Alexey Gritsenko, Xiao Wang,	920
	Muhammad Ferjad Naeem, Ibrahim Alabdulmohsin,	921
	Nikhil Parthasarathy, Talfan Evans, Lucas Beyer,	922
	Ye Xia, Basil Mustafa, and 1 others. 2025. Siglip	923
	2: Multilingual vision-language encoders with im-	924
	proved semantic understanding, localization, and	925
	dense features. <i>arXiv preprint arXiv:2502.14786</i> .	926
	Daniel Yang, Yao-Hung Hubert Tsai, and Makoto Ya-	927
	mada. 2024. On verbalized confidence scores for	928
	llms. <i>arXiv preprint arXiv:2412.14737</i> .	929
	Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri,	930
	Dan Jurafsky, and James Zou. When and why vision-	931
	language models behave like bags-of-words, and	932
	what to do about it? In <i>The Eleventh International</i>	933
	<i>Conference on Learning Representations</i> .	934
	Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu,	935
	Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan,	936
	Weijie Su, Jie Shao, and 1 others. 2025. InternV3:	937
	Exploring advanced training and test-time recipes	938
	for open-source multimodal models. <i>arXiv preprint</i>	939
	<i>arXiv:2504.10479</i> .	940
	A Visual De-Noising System Prompt	941
	To ensure reproducibility of the Symbolic Anchors	942
	(<i>Vidiodom</i>), we provide the exact system instructions	943
	used to transform the SemEval-2025 dataset im-	944
	ages.	945
	Task: Analyze the input image and trans-	946
	form it into a minimalist, abstract sym-	947
	bolic icon.	948
	1. Conceptual Instructions (De-	949
	Noising):	950
	• Identify the Core Essence: Deter-	951
	mine the fundamental meaning or	952
	action of the image. Ignore spe-	953
	cific details, individuals, or environ-	954
	ments.	955
	• Abstract & Merge (Metonymy):	956
	If the image contains multiple el-	957
	ements forming a narrative, distill	958
	them into a single, unified glyph	959
	that represents the entire concept	960
	(e.g., instead of “person watching	961
	loud TV,” create a symbol for “in-	962
	tense viewing”).	963
	• Remove Context: Eliminate all	964
	background elements, environ-	965
	ments, and secondary objects.	966
	2. Stylistic Instructions (Flat Iconog-	967
	raphy):	968

- 969 • **Geometric Reconstruction:** Re-
970 build the concept using only pure
971 geometric primitives (perfect cir-
972 cles, squares, triangles, and clean,
973 uniform arcs). Avoid organic or
974 sketchy lines.
- 975 • **Strict Flat Design:** There must be
976 absolutely zero gradients, shadows,
977 textures, or lighting effects. All col-
978 ors must be solid flats.
- 979 • **Bold Outlines:** Encase all major el-
980 ements in thick, uniform black out-
981 lines.
- 982 • **Limited Palette:** Restrict the color
983 palette strictly to Black, White, and
984 a maximum of two highly contrast-
985 ing solid accent colors derived from
986 the most prominent color in the in-
987 put image.
- 988 • **Composition:** The final output
989 should be a clean, centered logo
990 icon on a plain white background.

991 **B Examples**

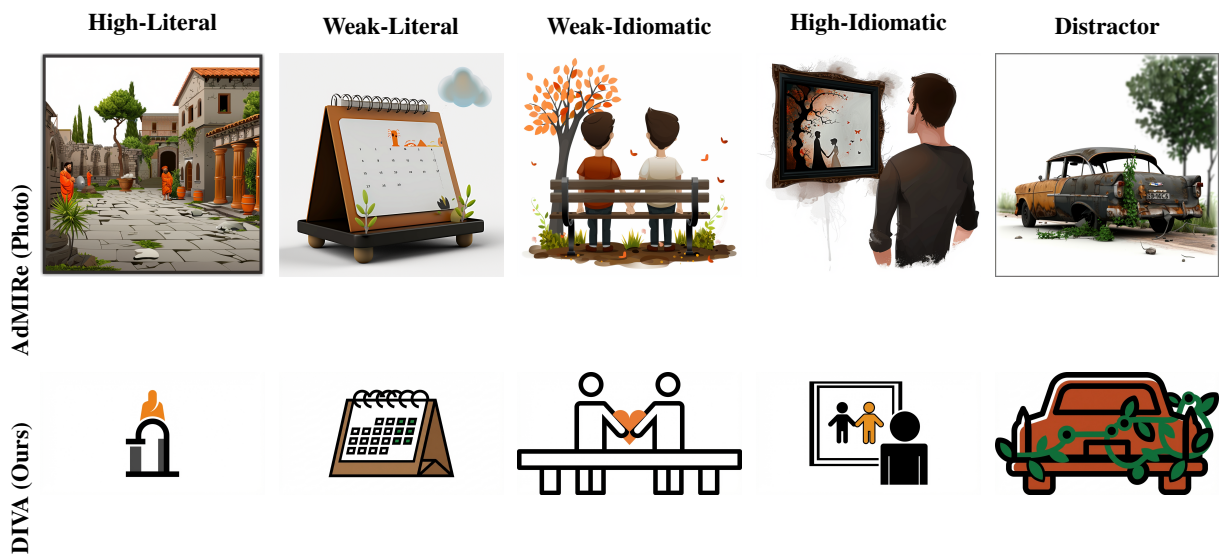


Figure 4: **Visual De-Noising in Action (AdMIRe vs. DIVA)**. Top Row: The original photorealistic images from ADMIRE, where high-frequency texture creates “semiotic noise.” Bottom Row: Our corresponding DIVA icons. By systematically de-noising the images across the full semantic spectrum (from Literal to Idiomatic), DIVA provides a clean, structure-aware testbed for multimodal reasoning.