

TRANSBIND: Explainable Compositional Grounding in LVLMs via Agentic Counterfactual Diagnostics

Anonymous ACL submission

Abstract

Large vision–language models (LVLMs) produce fluent but often unfaithful outputs—correct objects with wrong attributes, reversed relations, or mis-scoped negations. We frame these as *compositional grounding failures* and present TRANSBIND, a framework making such errors testable via *counterfactual minimal pairs* that change one semantic factor while holding others fixed. We introduce two metrics: SEC (selective edit consistency) and AFC (attribution factor coherence), plus a lightweight structured binding module. Experiments on four LVLMs show consistent improvements.

1 Introduction

LVLMs built on vision–language pretraining (Radford et al., 2021; ?; Li et al., 2023a) and instruction tuning (Liu et al., 2023a; Dai et al., 2023a) achieve strong performance but suffer from hallucinations (Rohrbach et al., 2018; Li et al., 2023b; Dai et al., 2023b). Many errors are *component-level*: the object is correct but the color is wrong; entities are correct but relations reversed; negations are ignored. We view these as *binding failures*—the model fails to bind semantic roles (attributes, relations, operators) to correct visual evidence.

This motivates testing *stability under controlled interventions*: when a minimal pair changes one factor, a grounded model should update only the corresponding output, and explanations should shift to relevant evidence.

Contributions. We propose TRANSBINDwith: (1) DIAGTRANS, an agentic diagnostic suite generating counterfactual minimal pairs targeting attributes, relations, quantities, and negation; (2) SEC/AFC metrics for output consistency and explanation coherence; (3) a structured binding module using object-centric slots (Locatello et al., 2020) and factorized cross-attention.

Factor	Pattern and expected behavior
Attribute	Flip mug color; only color token should change; attribution stays on mug.
Relation	Flip <i>on</i> → <i>under</i> ; truth value flips; referents stable.
Negation	Change <i>two</i> → <i>three</i> or add <i>not</i> ; only operator changes.

Table 1: Minimal-pair patterns: selective stability under single-factor interventions.

2 Failure Taxonomy and Minimal Pairs

Based on prior work (Rohrbach et al., 2018; Li et al., 2023b; Dai et al., 2023b; Thrush et al., 2022), we organize grounding failures into: (1) **entity failure**—mentioning unsupported objects; (2) **attribute binding failure**—correct object, wrong attribute; (3) **relational binding failure**—correct entities, wrong relation; (4) **quantifier failure**—wrong counts; (5) **logical operator failure**—ignored negation. Object-level metrics (CHAIR, POPE) capture (1); DIAGTRANStargets (2)–(5) via minimal pairs.

Table 1 illustrates minimal-pair patterns: each intervention changes one factor while expecting selective output change.

3 Related Work

Hallucination and compositionality. CHAIR (Rohrbach et al., 2018) and POPE (Li et al., 2023b) quantify object hallucination. Compositional benchmarks (Winoground (Thrush et al., 2022), VALSE (Parcalabescu et al., 2022), NLVR2 (Suhr et al., 2019), CLEVR (Johnson et al., 2017)) probe systematic grounding but lack counterfactual structure.

Explainability. Attribution methods (Sundararajan et al., 2017; Selvaraju et al., 2016) visualize evidence, but attention may not explain (Jain and Wallace, 2019). ERASER (DeYoung et al., 2020)

Pipeline: Seed \rightarrow Agent proposes T_f \rightarrow Instantiate pair \rightarrow Verify \rightarrow Evaluate with SEC/AFC

Figure 1: TRANSBINDpipeline overview.

highlights the gap between plausible rationales and faithful evidence. TRANSBINdevaluates explanations via counterfactual tests.

4 Problem Setup

Let $p_\theta(y | I, x)$ be an LVLM where I is an image and x is a question/instruction. We view output y as a sequence of semantic components (object mentions, attributes, relations, operators).

Counterfactual minimal pairs. An intervention T_f changes factor f while keeping others fixed:

$$(I, x) \text{ and } T_f(I, x) = (I', x').$$

In practice, I' can be an edited image, a paired image from a dataset, or the same image with a counterfactual prompt.

Selective stability principle. A grounded model should satisfy: (1) **targeted update**— f -related output changes appropriately; (2) **invariance elsewhere**—unrelated components remain stable; (3) **evidence shift**—attribution shifts to f -relevant regions.

5 DIAGTRANS: Counterfactual Diagnostics

DIAGTRANS targets four factor families: attributes (color/material), relations (spatial), quantification, and negation. Data sources include COCO (Lin et al., 2014), Visual Genome (Krishna et al., 2017), GQA (Hudson and Manning, 2019), and NLVR2 (Suhr et al., 2019).

We support two modes: **paired-image** (mining images differing in one factor) and **paired-prompt** (flipping operators while keeping the image fixed). An LLM agent proposes candidates following templates; verification filters non-minimal pairs using dataset annotations, lightweight detectors (Liu et al., 2023b; Kirillov et al., 2023), and human audit. Each pair carries metadata: target factor, confound checklist, verification evidence, and minimality score.

6 Metrics

We introduce two metrics for counterfactual evaluation over minimal pairs.

6.1 SEC: Selective Edit Consistency

Given outputs y, y' , we compute edit alignment via Levenshtein backtrace. Let $\Delta(y, y')$ be edited tokens and S_f be factor- f tokens (extracted via lexicon + pattern rules). We define:

$$\text{SEC} = \frac{|\Delta(y, y') \cap S_f|}{|\Delta(y, y')| + \epsilon},$$

where ϵ avoids division by zero. High SEC indicates the model changes only what should change. We canonicalize factor mentions (e.g., “navy” \rightarrow “blue”) before extraction.

6.2 AFC: Attribution Factor Coherence

For attribution method A (e.g., integrated gradients, Grad-CAM), let $A(I, x)$ produce an attribution map. AFC measures whether attribution shifts to f -relevant regions R_f :

$$\text{AFC} = \text{sim}(A_f, R_f) - \lambda \text{sim}(A_{-f}, R_f),$$

where A_f and A_{-f} are attributions for factor-related and unrelated tokens. This makes explanation evaluation testable rather than anecdotal.

7 Structured Binding Module

Diagnostics are useful only if they inform improvements. We describe a modular inductive bias for explicit role-filler binding.

7.1 Motivation

Classical compositional views treat meaning as role-filler binding (Smolensky, 1990). In LVLMs, the challenge is binding linguistic roles (“color of the mug”) to correct visual referents. Generic cross-attention can represent this implicitly but does not guarantee stable binding under interventions.

7.2 Design

We extract object-centric slots $S = \{s_i\}$ via Slot Attention (Locatello et al., 2020), then apply *factorized cross-attention*: object, attribute, and relation queries attend separately to slots before composition (Algorithm 1). This reduces the chance that frequent language patterns override weak visual signals.

7.3 Training Signals

DIAGTRANS can provide weak supervision via two optional objectives: (1) **selective invariance**—penalize changes in H_{-f} under T_f ; (2) **factor-aligned change**—encourage controlled change in H_f via contrastive margin.

Algorithm 1 Factorized Binding Fusion

Require: Image features V , text tokens X
 1: Slots $S \leftarrow \text{SlotAttention}(V)$
 2: $Q_{\text{obj}}, Q_{\text{attr}}, Q_{\text{rel}} \leftarrow \text{Proj}(X)$
 3: $H_i \leftarrow \text{CrossAttn}(Q_i, S)$ for $i \in \{\text{obj}, \text{attr}, \text{rel}\}$
 4: **return** $\text{Compose}(H_{\text{obj}}, H_{\text{attr}}, H_{\text{rel}})$

Family	Pass%	κ	Notes
Attribute	88	.78	color cues
Relation	81	.71	viewpoint
Quantifier	84	.74	counting
Negation	79	.69	scope

Table 2: Human audit ($n=600$).

8 Experiments

8.1 Benchmark Instantiation

We instantiate DIAGTRANS-ACL from COCO (Lin et al., 2014), GQA (Hudson and Manning, 2019), NLVR2 (Suhr et al., 2019), VALSE (Parcalabescu et al., 2022), and Winoground (Thrush et al., 2022). The suite contains 8,400 single-factor pairs. Figure 2 shows the composition. Human audit on 600 pairs (Table 2) achieves 76–88% minimality pass rate with Cohen’s κ 0.67–0.78.

8.2 Models

We evaluate four representative open LVLMs:

- **BLIP-2** (Li et al., 2023a): frozen vision encoder + Q-Former.
- **InstructBLIP** (Dai et al., 2023a): instruction-tuned BLIP-2.
- **LLaVA** (Liu et al., 2023a): visual instruction tuning.
- **MiniGPT-4** (Zhu et al., 2023): frozen encoder + LLM alignment.

8.3 Methods Compared

We compare: (1) **Base**—original LVLM; (2) **+TRANSBIND**—add factorized binding fusion; (3) **+TRANSBIND+cf-reg**—additionally apply counterfactual regularization.

8.4 Main Results

Table 3 shows TRANSBIND improves SEC and AFC across all models and factors, with largest gains on relations and negation.

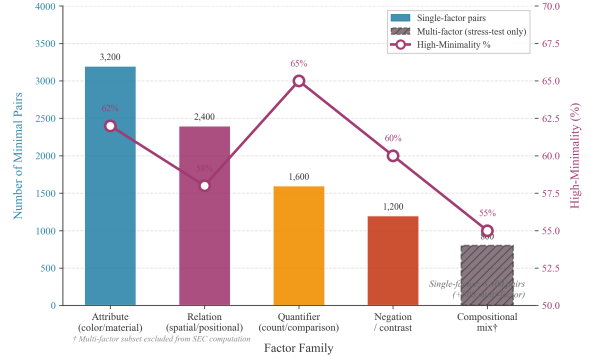


Figure 2: DIAGTRANS-ACL composition by factor family.

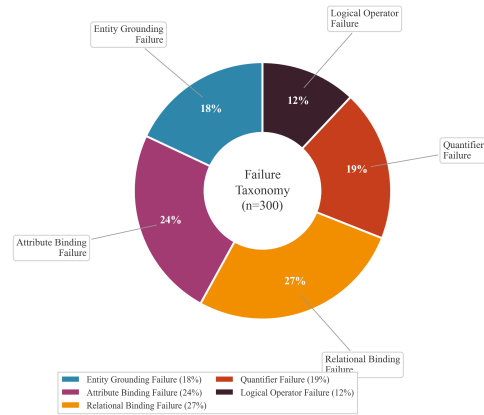


Figure 3: Error distribution on 300 audited pairs.

8.5 Ablations and Analysis

Table 4 shows rankings are stable across minimality levels. Table 5 confirms both slots and factorization contribute. SEC correlates with POPE ($\rho = -.62$) and CHAIR ($\rho = -.54$), suggesting compositional grounding improvements reduce hallucinations. Figure 3 shows binding failures dominate over entity hallucination.

8.6 Implementation Details

We use $K=7$ slots with 3 refinement iterations. Factor queries are projected via a 2-layer MLP (GELU, hidden 512). Counterfactual regularization fine-tuning uses learning rate $1e-5$, batch size 64, and 1 epoch on $\sim 27k$ augmented instances.

Limitations

First, minimality is hard to guarantee in the wild: real images correlate attributes and relations, and editing images at scale is non-trivial. We mitigate this via a high-minimality subset and metadata-driven analysis. Second, SEC depends on factor token identification, which can be imperfect for

Model	SEC \uparrow				AFC \uparrow	
	Attr.	Rel.	Quant.	Neg.	IG	GradCAM
BLIP-2	.612	.438	.524	.387	.281	.253
+ TRANSBIND	.658	.527	.571	.476	.342	.308
+ TRANSBIND+ cf-reg	.687	.564	.598	.518	.371	.339
InstructBLIP	.641	.472	.553	.421	.304	.271
+ TRANSBIND	.683	.551	.592	.503	.358	.327
+ TRANSBIND+ cf-reg	.709	.589	.621	.541	.392	.358
LLaVA	.623	.451	.508	.406	.268	.241
+ TRANSBIND	.671	.538	.561	.489	.331	.297
+ TRANSBIND+ cf-reg	.698	.576	.593	.527	.364	.331
MiniGPT-4	.598	.426	.497	.378	.259	.232
+ TRANSBIND	.647	.516	.548	.467	.318	.286
+ TRANSBIND+ cf-reg	.678	.557	.581	.509	.352	.321

Table 3: Main results on DIAGTRANS-ACL (8,400 pairs).

Method	SEC (full)	SEC (high-min.)
Base	.508	.561
+ TRANSBIND	.559	.608
+ TRANSBIND+ cf-reg	.591	.642

Table 4: Robustness: full vs. high-minimality subset.

Variant	SEC	AFC
Base	.522	.288
+ slots only	.563	.324
+ factorization only	.571	.331
+ TRANSBIND(both)	.592	.358

Table 5: Ablation on InstructBLIP.

paraphrased outputs; we use canonicalization but edge cases remain. Third, AFC requires region proxies (annotations or detectors), which may inject their own biases. These limitations motivate transparent reporting and ablation of verification strength.

Ethical Considerations

More faithful LVLMs can improve safety in assistive and decision-support settings, but diagnostic tools can also be misused to create adversarial prompts. We recommend releasing diagnostics with clear intended use and avoiding sensitive images or personal data. When human verification is used, annotators should be fairly compensated and protected from harmful content exposure.

9 Conclusion

We presented TRANSBIND, a framework for making compositional grounding in LVLMs measurable and explainable through agentic counterfac-

tual diagnostics. By linking minimal interventions to output consistency (SEC) and explanation coherence (AFC), TRANSBIND provides an evaluation loop that can guide architectural improvements. Experiments show consistent gains across four LVLm families.

Acknowledgments

We thank the authors of the public datasets and benchmarks used in this work.

References

- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven C. H. Hoi. 2023a. [Instructblip: Towards general-purpose vision-language models with instruction tuning](#). *CoRR*, abs/2305.06500.
- Wenliang Dai, Zihan Liu, Ziwei Ji, Dan Su, and Pascale Fung. 2023b. [Plausible may not be faithful: Probing object hallucination in vision-language pre-training](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2136–2148, Dubrovnik, Croatia. Association for Computational Linguistics.
- Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani, Eric Lehman, Caiming Xiong, Richard Socher, and Byron C. Wallace. 2020. [ERASER: A benchmark to evaluate rationalized NLP models](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 4443–4458. Association for Computational Linguistics.
- Drew A. Hudson and Christopher D. Manning. 2019. [GQA: A new dataset for real-world visual reasoning and compositional question answering](#). In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20,*

258	2019, pages 6700–6709. Computer Vision Foundation / IEEE.	315
259		316
260	Sarthak Jain and Byron C. Wallace. 2019. Attention is not explanation . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)</i> , pages 3543–3556. Association for Computational Linguistics.	317
261		318
262		319
263		320
264		321
265		322
266		323
267		324
268		325
269	Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross B. Girshick. 2017. CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning . In <i>2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017</i> , pages 1988–1997. IEEE Computer Society.	326
270		327
271		328
272		329
273		330
274		331
275		332
276		333
277	Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloé Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross B. Girshick. 2023. Segment anything . <i>CoRR</i> , abs/2304.02643.	334
278		335
279		336
280		337
281		338
282		339
283	Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations . <i>Int. J. Comput. Vis.</i> , 123(1):32–73.	340
284		341
285		342
286		343
287		344
288		345
289	Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. 2023a. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models . <i>CoRR</i> , abs/2301.12597.	346
290		347
291		348
292		349
293		350
294		351
295		352
296		353
297		354
298		355
299		356
300		357
301		358
302		359
303		360
304		361
305		362
306		363
307		364
308		365
309		366
310		367
311		368
312		369
313		370
314		371