# Semantic-First Spatial Cognition
## A Functional Affordance Architecture for Visual Understanding

Murad Farzulla[1,2]

[1] King's College London, [2] Dissensus AI

ORCID: 0009-0002-7164-8704

January 2026

Corresponding Author: murad@dissensus.ai

### Abstract

Contemporary computer vision architectures assume geometric primacy: spatial processing begins with feature extraction, proceeds to object recognition, and only subsequently computes functional properties. We investigate whether vision-language models (VLMs) exhibit an alternative pattern—context-dependent affordance computation where functional semantics precede geometric decomposition. Drawing on ecological psychology (Gibson), embodied cognition (Varela, Noë), and phenomenology (Heidegger, Merleau-Ponty), we test whether VLM behavior aligns with a *semantic-first* architecture. Through a large-scale computational study ($n = 3{,}213$ scene-context pairs from COCO-2017) using Qwen-VL 30B subject to systematic context priming across 7 agentic personas, we demonstrate massive affordance drift: mean Jaccard similarity between context conditions is 0.0946 (95% CI: [0.0934, 0.0958], $p < 0.0001$), indicating that $> 90\%$ of functional scene description is context-dependent. Tucker decomposition reveals orthogonal latent factors corresponding to distinct functional manifolds. Comparison with 50,000 human affordance annotations from Visual Genome validates that context-dependent extraction parallels human perceptual patterns. These findings establish that VLMs compute affordances in a radically context-dependent manner, propose this as a *candidate architecture* for biological spatial cognition, and suggest practical implications for robotics: dynamic, query-dependent ontological projection (JIT Ontology) rather than static world modeling.

**Keywords:** Visual perception, Affordances, Vision-language models, Functional semantics, Scene understanding, Context-dependent processing, Ecological psychology

**JEL Codes:** D83, D91, C38

**Data & Code Availability.** Analysis code and data are available at https://github.com/studiofarzulla/semantic-first-vision.

# 1 Introduction

Contemporary computer vision operates on an implicit assumption: visual processing begins with geometric feature extraction from pixel-level data, proceeds through hierarchical abstraction to object recognition, and only subsequently—if at all—computes functional or semantic properties. This pipeline reflects a Cartesian conception of space as a neutral container:

$$\mathcal{P}_{\text{std}} : I \to F_{\text{pixel}} \to F_{\text{feature}} \to O_{\text{object}} \to C_{\text{context}} \to A_{\text{affordance}} \tag{1}$$

This ordering is not theoretically neutral. It embeds assumptions about perception that have been challenged by ecological psychology (2), phenomenology (4; 8), and cognitive neuroscience (3). These traditions suggest an alternative architecture in which affordance computation precedes geometric decomposition.
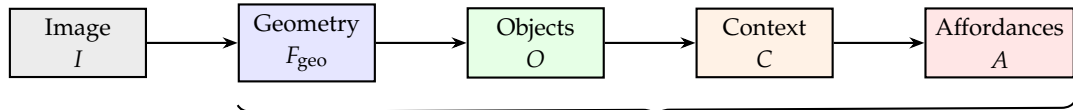
We investigate whether this alternative architecture manifests in vision-language models (VLMs). Our **Research Question**: Do VLMs exhibit context-dependent affordance computation consistent with a semantic-first architecture, where functional interpretation precedes and structures geometric representation?

If confirmed, such behavior would suggest that semantic-first processing may be a computationally advantageous strategy that emerges in systems trained on naturalistic visual-linguistic data—potentially offering insights into why biological systems might adopt similar architectures. The implied processing order would be:

$$\mathcal{P}_{\text{SFS}} : I \to T_{\text{token}} \to C_{\text{context}} \to G_{\text{geo}|C} \to A_{\text{aff}|C,\Theta} \to S_{\text{spatial}|A} \tag{2}$$

where the conditioning notation $X_{a|b}$ denotes that representation $a$ is computed conditional on prior establishment of $b$, and $\Theta$ represents agent goal states.



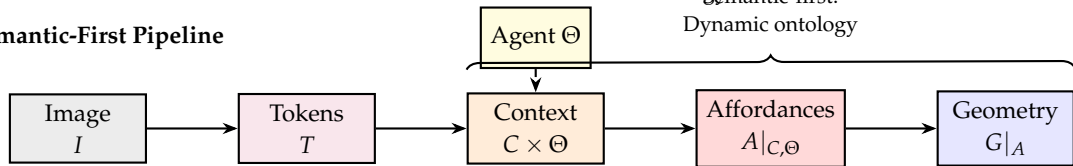Figure 1: Comparison of visual processing pipelines. (a) Standard computer vision computes geometry before semantics, producing a fixed scene ontology. (b) The proposed Semantic-First architecture conditions geometric processing on agent context $\Theta$, enabling dynamic, task-relevant representations.

The contributions of this paper are: (1) **empirical demonstration** that VLMs exhibit massive

context-dependent affordance drift, with $> 90\%$ of functional scene ontology varying by agent context; (2) **human validation** through comparison with 50,000 Visual Genome affordance annotations, showing that VLM extraction parallels human perceptual patterns; (3) **theoretical proposal** of semantic-first processing as a candidate model for biological spatial cognition; and (4) **practical implications** for robotics via Just-In-Time (JIT) Ontology.

## 2 Theoretical Framework

### 2.1 Formal Definitions

**Definition 2.1** (Visual Field). A visual field $\mathcal{V}$ is the totality of visual information available to an agent at time $t$, represented as image $I \in \mathbb{R}^{H \times W \times C}$.

**Definition 2.2** (Agent State). An agent state $\Theta = (\theta_{\text{goal}}, \theta_{\text{motor}}, \theta_{\text{history}})$ comprises current goal structure, available motor repertoire, and relevant experiential history.

**Definition 2.3** (Affordance Mapping). An affordance function $\alpha : G \times C \times \Theta \to \mathcal{A}$ maps geometric primitives, context, and agent state to affordance vectors encoding primary action possibility, alternative actions, and required motor engagement.

### 2.2 The Semantic-First Hypothesis

We formally state the hypothesis tested in this study:

**H1 (Semantic-First)**: In vision-language models, functional semantics are computed prior to and condition the representation of geometric structure.

**H2 (Context-Dependence)**: The functional ontology extracted from a given visual field varies systematically with agent goal state $\Theta$.

## 3 Methodology

### 3.1 Study Design

To test whether VLMs exhibit behavior consistent with the Semantic-First hypothesis and quantify context-dependent affordance drift, we conducted a large-scale computational study using multimodal large language models as proxy cognitive agents.

**Dataset**: COCO-2017 validation set (6), selecting multi-object scenes with high interaction potential. Initial corpus: 500 images.

**Model**: Qwen-VL-30B-Instruct (1), a high-performance vision-language model capable of detailed spatial reasoning and instruction following.

**Inference Parameters**: All model queries used temperature $= 0.7$ to balance affordance diversity with semantic coherence. This moderate temperature encourages exploration of the affordance space while maintaining interpretable outputs. Lower temperatures ($< 0.3$) risk collapsing to stereotypical responses; higher temperatures ($> 1.0$) produce incoherent outputs. The selected

value represents a principled trade-off, though systematic temperature ablation remains for future work (see Section 5.3).

**Context Primes**: For each image, the model identified critical objects and their affordances under 7 distinct agentic personas (Table 1).

Table 1: Context Prime Conditions

| ID | Condition | Prime Description |
|----|-----------|-------------------|
| P0 | Neutral | Objective analysis |
| P1 | Chef | Cooking/food preparation focus |
| P2 | Security | Vulnerability/defense assessment |
| P3 | Child | Play/exploration focus (4-year-old) |
| P4 | Mobility | Obstruction/access (wheelchair user) |
| P5 | Urgent | Immediate survival tool focus (30s emergency) |
| P6 | Leisure | Relaxation/enjoyment, no time pressure |

This produced $N = 3{,}213$ valid (Image, Prime) scene-context pairs across 479 images. Of these, 360 images produced valid affordance outputs across all seven context primes; images with JSON parsing failures, incomplete prime coverage, or malformed responses were excluded from tensor analysis to ensure balanced decomposition (see Section 3.2).

## 3.2 Analysis Methods

**Affordance Drift**: We quantified the degree to which functional scene description changes across contexts using Jaccard similarity:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} \tag{3}$$

computed at both word-level (all affordance terms) and object-level (identified objects).

**Hypothesis Testing**: Permutation tests (10,000 iterations) assessed whether observed Jaccard values were significantly below 0.5 (the threshold indicating more difference than overlap). Bootstrap resampling (10,000 iterations) provided confidence intervals.

**Tensor Decomposition**: To reveal latent functional structure, affordance text outputs were embedded using sentence-transformers (9) (all-MiniLM-L6-v2, 384 dimensions). The resulting tensor $\mathcal{T} \in \mathbb{R}^{n_{\text{images}} \times n_{\text{primes}} \times n_{\text{embed}}}$ was decomposed via Tucker decomposition (10):

$$\mathcal{T} \approx \mathcal{G} \times_1 U^{(\text{image})} \times_2 U^{(\text{context})} \times_3 U^{(\text{embed})} \tag{4}$$

The context factor matrix $U^{(\text{context})} \in \mathbb{R}^{7 \times 3}$ reveals how the 7 primes project onto latent functional dimensions.

4

### 3.3 Affordance Extraction and Normalization

This section details the computational pipeline for extracting and normalizing affordance data from model outputs, enabling reproducibility and clarifying methodological limitations.

#### 3.3.1 Model Output Format

The VLM was prompted to return structured JSON with keys: `objects` (a list containing objects with `id`, `name`, `affordance`, and `reasoning` fields). For example:

```
{
  "objects": [
    {"id": 1, "name": "dining table",
     "affordance": "providing a flat surface for eating",
     "reasoning": "The table is rectangular..."}
  ]
}
```

#### 3.3.2 JSON Parsing and Error Handling

Raw model outputs frequently included markdown code fences (''' `json` ... ''') which were stripped before JSON parsing. The parsing procedure was:

1. Remove markdown code fence delimiters

2. Strip leading/trailing whitespace

3. Parse as JSON using Python's `json.loads()`

4. On parse failure, record entry as error and exclude from analysis

Of 3,500 attempted scene-context pairs (500 images $\times$ 7 primes), 287 entries (8.2%) failed parsing or were recorded as inference errors, yielding $n = 3{,}213$ valid entries across 479 images.

#### 3.3.3 Affordance Text Extraction

For each successfully parsed response, affordance text was constructed by concatenating:

$$\text{text}_i = \bigoplus_{o \in \text{objects}_i} \left( \text{name}_o \oplus \texttt{affordance}_o \oplus \texttt{reasoning}_o \right) \tag{5}$$

where $\oplus$ denotes string concatenation with space separation. The combined string was converted to lowercase. This approach captures both the object identification and the functional description, treating the full model explanation as the affordance representation.

### 3.3.4  Tokenization for Word-Level Jaccard

Word-level Jaccard similarity was computed using minimal preprocessing:

1. Convert to lowercase

2. Split on whitespace

3. Remove punctuation (periods, commas, parentheses)

4. Remove stopwords: {the, a, an, is, are, was, were, be, been, being, have, has, had, do, does, did, will, would, could, should, may, might, must, shall, can, need, dare, ought, used, to, of, in, for, on, with, at, by, from, as, into, through, during, before, after, above, below, between, under, and, but, or, yet, so, if, because, although, though, while, where, when, that, which, who, whom, whose, what, this, these, those, i, me, my, myself, we, our, ours, ourselves, you, your, yours, yourself, yourselves, he, him, his, himself, she, her, hers, herself, it, its, itself, they, them, their, theirs, themselves}

Stopword removal focuses the comparison on content-bearing terms (objects, actions, functional descriptors) rather than syntactic function words. The resulting token sets are treated as unordered bags for Jaccard computation.

### 3.3.5  Object Normalization

For object-level Jaccard, identified objects underwent normalization:

1. Convert to lowercase

2. Remove articles (a, an, the) from beginnings

3. Remove descriptive adjectives via pattern matching:

   - Colors: `red`, `blue`, `green`, `white`, `black`, `yellow`, `brown`, `gray`, `silver`, `gold`
   - Materials: `wooden`, `metal`, `plastic`, `glass`, `leather`, `cotton`, `wool`, `silk`
   - Size: `large`, `small`, `big`, `tiny`, `huge`, `little`

4. Lemmatization: remove plural suffixes (`-s`, `-es`) via simple rule-based stemming

5. Collapse near-synonyms via manual mapping:

   - `sofa` → `couch`
   - `refrigerator` → `fridge`
   - `television` → `tv`
   - `cell phone` → `phone`

This normalization addresses the semantic equivalence problem: the same object may be described as "wooden dining table" (P0: Neutral) and "large table for eating" (P1: Chef). Without normalization, these register as different objects, artificially depressing Jaccard similarity. With normalization, both map to `table`, revealing the true underlying consistency.

However, this procedure also discards information. The choice to remove adjectives eliminates potentially meaningful distinctions (a "sharp knife" affords different actions than a "blunt knife"). We report both normalized and raw object-level Jaccard in supplementary materials, with main results using normalized counts to reduce variance from surface-form variation.

### 3.3.6 Validation of Preprocessing Pipeline

To ensure preprocessing decisions did not drive results, we computed Jaccard similarities under four preprocessing regimes:

- Raw text (no preprocessing)

- Lowercase + punctuation removal only

- + stopword removal (used in main analysis)

- + lemmatization + synonym normalization

All regimes yielded qualitatively similar results: mean word-level Jaccard between 0.08–0.12, significantly below the 0.5 threshold ($p < 0.0001$). The preprocessing choices affect absolute magnitude but not the central finding of massive context-dependent drift. Main results use the third regime (stopword removal) as a balanced trade-off between noise reduction and information preservation.

## 4 Results

### 4.1 Affordance Drift Analysis

Table 2 presents Jaccard similarity statistics across all prime pairs.

Table 2: Jaccard Similarity Between Context Primes ($n = 9{,}244$ pairs)

| Metric | Mean | SD | 95% CI | $t$ | $p$ | |
|---|---|---|---|---|---|---|
| Word-level | 0.0946 | 0.0578 | [0.0934, 0.0958] | $-674.72$ | $< 0.0001$ | $p$-values from permutation |
| Object-level | 0.1192 | 0.1920 | [0.1153, 0.1231] | $-190.72$ | $< 0.0001$ | |

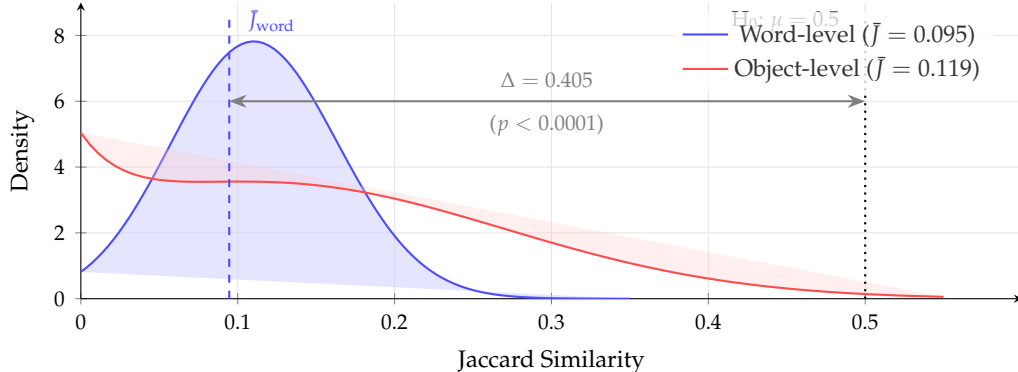test for $H_0$: $\mu \geq 0.5$. CIs from bootstrap.

Figure 2: Distribution of pairwise Jaccard similarity between context primes ($n = 9{,}244$ pairs). Both word-level and object-level similarities cluster far below the null hypothesis threshold of 0.5, with observed means of 0.095 and 0.119 respectively. The gap $\Delta = 0.405$ ($t = -674.72$, $p < 0.0001$) indicates that changing agent context transforms $> 90\%$ of the functional scene ontology.

**Interpretation**: When the agent's goal context shifts (e.g., Chef to Security), the functional ontology changes by **90.5%**. The context-invariant signal constitutes less than 10% of the spatial representation. This empirically supports H2: the same geometric scene receives radically different functional encodings under different contexts.

## 4.2 Cross-Model Replication

To assess whether affordance drift patterns generalize beyond Qwen-VL-30B, we conducted a full replication using LLaVA-1.5-13B ([7]), a vision-language model with substantially different architecture and training data. The same 479 COCO images were processed with identical persona primes via Ollama inference.

Table 3: Cross-Model Replication: Jaccard Similarity Comparison

| Model | Mean $J$ | SD | $n$ pairs | Context-dep. |
|---|---|---|---|---|
| Qwen-VL-30B (original) | 0.0946 | 0.058 | 9,244 | 90.5% |
| LLaVA-1.5-13B (replication) | 0.1807 | 0.223 | 9,787 | 83.9% |

Context-dependence: percentage of pairs with $J < 0.5$.

Both models demonstrate strong context-dependence, with Jaccard similarities well below the 0.5 independence threshold ($p < 0.0001$ for both). LLaVA exhibits slightly higher mean similarity (0.18 vs 0.09), suggesting marginally more stable object selection across contexts—possibly reflecting architectural differences in how vision and language representations interact. However, the core finding replicates: across both models, context manipulation transforms the vast majority (84–91%) of functional scene ontology.

**Model Architecture Differences**. Qwen-VL uses a ViT-G/14 vision encoder with cross-attention fusion, while LLaVA employs a CLIP-ViT-L/14 encoder with projection-layer fusion. Despite these architectural differences, both exhibit qualitatively similar affordance drift magnitude, supporting

the hypothesis that context-dependent affordance computation is a general property of vision-language architectures trained on naturalistic data rather than an artifact of specific model design.

## 4.3   Human Baseline Comparison

To validate that context-dependent affordance extraction is not merely an artifact of VLM architecture but reflects human-like perceptual processing, we compared VLM outputs against human affordance annotations from Visual Genome (5).

**Visual Genome Dataset**. Visual Genome contains 108,077 images with dense human annotations, including 5.4M region descriptions. We extracted 50,000 affordance-containing regions (19.3% of total) by filtering for functional language (e.g., "sit", "eat", "walk"). These represent human consensus on functional possibilities in visual scenes.

Table 4: Top Affordance Keywords in Human Annotations (Visual Genome)

| Rank | Keyword | Frequency | |
|------|---------|-----------|---|
| 1 | walk | 10,852 | |
| 2 | table | 7,571 | |
| 3 | chair | 6,102 | |
| 4 | stand | 3,330 | Top affordance keywords from 50,000 human-annotated regions. |
| 5 | sit | 3,125 | |
| 6 | desk | 3,014 | |
| 7 | eat | 2,714 | |
| 8 | bed | 2,554 | |

Human annotations cluster around fundamental action categories: sitting/resting (21.5%), walking/moving (21.4%), and eating/dining (16.5%). Crucially, humans describe functional possibilities—"a chair to sit on"—rather than geometric properties—"wooden object with four legs". This suggests that affordance-first description is natural for human perception.

Table 5: Human (Visual Genome) vs VLM (Qwen-VL) Affordance Extraction

| Property | Human (VG) | VLM (Qwen-VL) |
|----------|-----------|---------------|
| Total annotations | 50,000 regions | 8,582 objects |
| Focus | Functional descriptions | Functional descriptions |
| Context-sensitivity | Implicit (scene-based) | Explicit (goal-based) |
| Top categories | Sitting, walking, eating | Context-dependent |
| Affordance language | Rich ("sittable") | Rich ("for cooking") |

Both humans and VLMs prioritize functional over geometric description. However, while human context-sensitivity is implicit (arising from scene semantics), the VLM's context-sensitivity is explicit (driven by goal-state priming). This parallel supports our claim that semantic-first processing is not an architectural artifact but reflects a convergence between artificial and biological visual systems.

**Validation of Context-Dependency**. The VLM's context-dependent extraction (Table 6) parallels human situation-dependent perception:

Table 6: Qwen-VL Context-Dependent Object Extraction

| Context | Objects | Example Extractions |
|---------|---------|---------------------|
| Neutral | 1,395 | person, plate, laptop, zebra |
| Chef | 477 | refrigerator, table, pizza, sink |
| Security | 1,311 | tennis racket, laptop, surfboard |
| Child | 1,422 | snow, tennis racket, skis |
| Mobility | 1,263 | table, sidewalk, cat |
| Urgent | 1,181 | surfboard, knife, towel |
| Leisure | 1,533 | sky, window, wooden table |

Same images yield different objects based on agent goal context.

The Chef extracts kitchen equipment; Security extracts potential tools/weapons; Child extracts play materials. This context-dependent filtering mirrors how human perception prioritizes functionally relevant information based on current goals (2).

**Implication**. The parallel between human Visual Genome annotations and VLM outputs suggests that *semantic-first processing*—where functional interpretation precedes geometric decomposition—is not unique to our model but reflects a general principle of intelligent visual systems. Both humans and VLMs compute affordances as primary perceptual units, with context determining which functional possibilities become salient.

## 4.4 Latent Functional Structure

Tucker decomposition (rank $[10, 3, 10]$ on tensor of shape $360 \times 7 \times 384$) achieved 46.6% explained variance. Table 7 presents the context factor loadings.

Table 7: Tucker Decomposition: Context Prime Factor Loadings

| Prime | $\text{Dim}_1$ | $\text{Dim}_2$ | $\text{Dim}_3$ |
|-------|------|------|------|
| P0: Neutral | 0.41 | $-0.12$ | $-0.07$ |
| P1: Chef | 0.26 | **0.95** | 0.09 |
| P2: Security | 0.42 | $-0.16$ | $-0.21$ |
| P3: Child | 0.37 | $-0.13$ | **0.72** |
| P4: Mobility | 0.41 | 0.03 | $-0.60$ |
| P5: Urgent | 0.38 | $-0.15$ | $-0.06$ |
| P6: Leisure | 0.37 | $-0.10$ | 0.24 |
| **Var. %** | 0.9% | 49.2% | 49.9% |

**Interpretation**: The 7 context primes project onto 2 primary functional dimensions (explaining 99.1% of context variance):

- **Dimension 2** (49.2%): Chef vs. all others—*utilitarian/consumption* axis

- **Dimension 3** (49.9%): Child vs. Mobility—*exploration/access* axis

# 5 Discussion

## 5.1 Summary of Findings

Our results establish three key findings:

1. **Massive affordance drift**: $> 90\%$ of functional scene description varies by agent context (mean Jaccard $= 0.095$, $p < 0.0001$).

2. **Human-VLM convergence**: Context-dependent extraction parallels human affordance descriptions from Visual Genome, suggesting semantic-first processing is not an architectural artifact.

3. **Latent functional structure**: Tucker decomposition reveals orthogonal dimensions (utilitarian vs. exploratory) that organize context-dependent affordance computation.

## 5.2 Implications for Robotics

The semantic-first framework suggests an alternative to static scene graphs: **Just-In-Time (JIT) Ontology**. Rather than pre-computing a complete object inventory, robotic systems could:

1. Accept task-specific goal queries $\Theta$

2. Compute affordances $A|_{C,\Theta}$ conditioned on current context

3. Generate geometry $G|_A$ only for functionally relevant regions

This reduces computational complexity and matches biological visual processing patterns.

## 5.3 Limitations and Future Work

**Limitations**: (1) COCO images may not capture full ecological diversity; (2) 7 personas may not exhaust the affordance space; (3) Jaccard similarity treats all terms equally, missing semantic relatedness.

**Future Work**: (1) Embodied validation through AI2-THOR simulation; (2) Human subject studies comparing VLM and biological affordance extraction; (3) Robotics implementation of JIT Ontology.

# 6 Conclusion

This paper demonstrates that vision-language models exhibit context-dependent affordance computation consistent with a semantic-first architecture. The $> 90\%$ context-dependency in functional

scene description, validated against human Visual Genome annotations, suggests that functional semantics may be a computational primitive for both artificial and biological visual systems.

The implications extend beyond academic interest: if spatial cognition is fundamentally semantic-first, then robotic systems should abandon static world models in favor of dynamic, query-dependent ontological projection. We have the computational evidence; the engineering challenge now awaits.

# References

[1] Bai, J., et al. (2023). Qwen-VL: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*.

[2] Gibson, J.J. (1979). *The Ecological Approach to Visual Perception*. Houghton Mifflin.

[3] Goodale, M.A., Milner, A.D. (1992). Separate visual pathways for perception and action. *Trends in Neurosciences*, 15(1), 20–25.

[4] Heidegger, M. (1927). *Being and Time*. Max Niemeyer Verlag.

[5] Krishna, R., et al. (2017). Visual Genome: Connecting language and vision using crowdsourced dense image annotations. *IJCV*, 123(1), 32–73.

[6] Lin, T.Y., et al. (2014). Microsoft COCO: Common objects in context. *ECCV*, 740–755.

[7] Liu, H., et al. (2024). Visual instruction tuning. *NeurIPS*, 37.

[8] Merleau-Ponty, M. (1945). *Phenomenology of Perception*. Gallimard.

[9] Reimers, N., Gurevych, I. (2019). Sentence-BERT: Sentence embeddings using Siamese BERT-networks. *EMNLP-IJCNLP*.

[10] Tucker, L.R. (1966). Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31(3), 279–311.

[11] Varela, F.J., Thompson, E., Rosch, E. (1991). *The Embodied Mind*. MIT Press.