# Is CLIP ideal? No. Can we fix it? Yes!

Raphi Kang     Yue Song     Gerogia Gkioxari     Pietro Perona
California Institute of Technology
Pasadena, CA

rkang@caltech.edu

## Abstract

*Contrastive Language-Image Pre-Training (CLIP) is a popular method for learning multimodal latent spaces with well-organized semantics. Despite its wide range of applications, CLIP's latent space is known to fail at handling complex visual-textual interactions. Recent works attempt to address its shortcomings with data-centric or algorithmic approaches. But what if the problem is more fundamental, and lies in the geometry of CLIP? Toward this end, we rigorously analyze CLIP's latent space properties, and prove that no CLIP-like joint embedding space exists which can correctly do any two of the following at the same time: 1. represent basic descriptions and image content, 2. represent attribute binding, 3. represent spatial location and relationships, 4. represent negation. Informed by this analysis, we propose Dense Cosine Similarity Maps (DCSMs) as a principled and interpretable scoring method for CLIP-like models, which solves the fundamental limitations of CLIP by retaining the semantic topology of the image patches and text tokens. This method improves upon the performance of classical CLIP-like joint encoder models on a wide array of benchmarks. We share our code and data here for reproducibility:* https://github.com/Raphoo/DCSM_Ideal_CLIP

## 1. Introduction

Contrastive Language-Image Pre-Training (CLIP) is a widely used method for pretraining vision language model (VLM) embeddings in downstream applications [46]. It jointly trains an image and text encoder, mapping both modalities into a shared latent space. In this space, cosine similarity between embeddings reflects semantic similarity between the text and image: a high value indicates a semantic match between text and image, while a low value suggests that they are unrelated.

CLIP is simple, computationally efficient, and CLIP-based zero-shot multimodal tasks such as image classification and text-based retrieval can be impressively accu-
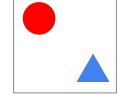


Figure 1. CLIP scores do not accurately reflect semantics of text prompts due to inherent geometric limitations. For five out of six pairs of captions, the incorrect pair has the higher CLIP score with the image. By contrast, our proposed model using Dense Cosine Similarity Maps (DCSM) correctly scores matched pairs. The similarity score is unnormalized because it is predicted by a neural network. CLIP scores computed with OpenAI-CLIP ViT-B/32.

rate. Unlike autoregressive models where an increasingly large number of images and texts to compare requires prohibitively longer context windows, CLIP can process each image and text prompt separately and only requires a simple inner product calculation for scoring. As a result, the list of systems that use CLIP embeddings toward more complex capabilities, such as image captioning [25], visual question answering [51], and text-guided image manipulation/generation [19, 38], gets longer by the day.

CLIP, however, is not perfect. It struggles with spatial reasoning and compositional understanding [18, 56, 66], concept and attribute binding [7, 22, 37], as well as negation [1, 52], to name a few shortcomings. Examples in Fig. 1 demonstrate these failure modes. These defects of CLIP impact downstream models and tasks. For example, search engines which use CLIP scores will bring up images of yellow coats to the prompt "not a yellow coat", and CLIP-based generative models cannot reflect spatial relationships in the text prompt. Ever since CLIP's conception, the vision-language community has worked to improve its semantics via adjustments in the training data [1, 20, 28, 40, 45, 52, 59, 60, 66] or the architecture

and training procedure [8, 26, 27, 35, 50, 58, 69]. But could we be pushing Sisyphus's boulder? In this work, we suggest to take a step back and reassess CLIP's basic geometry and philosophy from first principles.

We ask two questions: (1) *Is it possible for there to exist a CLIP-like latent space which has all the properties necessary to succeed in popular VLM tasks, using cosine similarity as the reference metric?* (2) *If no such space exists - is there a way to "rescue" this latent space without forsaking CLIP altogether?*

To address these queries, we first formalize the geometry of the CLIP latent space and establish a list of conditions which must be satisfied in order for it to understand the precise meaning of images and texts, as required by popular VLM benchmarks [1, 15, 18, 22, 56, 66]. Second, we show that these conditions are *not achievable* when using cosine similarity on the unit hypersphere. This answers the first question: an ideal CLIP space *cannot* exist. Finally, we propose a lightweight downstream CNN which utilizes a topological map relying on cosine similarity scores between each text token and image patch embedding, to produce a more perceptive text-image distance score in lieu of naive cosine similarity in CLIP-space. With this simple adjustment, we find our simple model can outperform state of the art joint-encoder models on a wide array of benchmarks.

In summary, our contributions are:

- **Problem Identification**: We find that naive cosine similarity on unit vector embeddings have fundamental geometric restrictions preventing it from accurately representing (1) attribute binding to distinct object concepts, (2) spatial relationships and location, and (3) negation.
- **Analysis/Proof**: We define the CLIP latent space as a projection of images and texts which are composed of atomic objects, attributes, and spatial relationships, and formalize desired conditions for the latent space in order for it to succeed on existing VLM benchmarks. We prove for each condition that no satisfactory vector space exists.
- **Topology as a Solution**: We propose to use a Dense Cosine Similarity Map (DCSM) from pre-trained CLIP to produce a more comprehensive text-image score. We present a prototype with a simple two-layer CNN module.
- **Experimental evaluation**: We evaluate our method on multiple benchmark regimes against CLIP-like models, and observe consistent performance gain across tasks.

## 2. Prior Work

**Vision Language Models.** Building upon the principles of large-scale contrastive pretraining and joint representation learning, recent CLIP-like VLMs [16, 24, 36, 64, 64, 65] have significantly bridged the gap between vision and language. Beyond CLIP models, autoregressive VLMs [31, 32] have emerged as compelling alternatives by jointly attending to both text and image embeddings. Neurosymbolic

program synthesis methods [14, 34, 54] also mitigate some of CLIP's limitations by formulating text-image semantic distance acquisition as several sub-problems. While these methods are more comprehensive than CLIP in complex visual reasoning, most of them rely on a latent space of CLIP-like models as a core submodule. As such, there remains a strong motivation to continue improving CLIP architectures by addressing their inherent shortcomings.

**Empirical Limitations of CLIP.** A growing body of work has revealed several limitations of CLIP in handling complex visual-text interactions. Specifically, CLIP struggles with accurately binding attributes to concepts in multi-object scenes [7, 22, 37], exhibits misinterpretations of object layouts or conflates multiple entities within a single scene [22, 37, 66], and has difficulty in accurately representing negation [1, 52]. To systematically evaluate these limitations, various vision language benchmarks have been proposed targeting different tasks, including attributes, relationships, and order (ARO) [66], Sugarcrepe [15], VL-checklist [72], WhatsUp [18], Multimodal Visual Patterns (MMVP) [55], and NegBench [1].

**Geometric Analysis of CLIP.** Several studies have examined the geometric properties of CLIP. Some works empirically assess the CLIP latent space and find that existing CLIP models underutilize the latent space by having a modality gap between text and image embeddings [29], being highly anisotropic even within the same modality [21], and that CLIP training is empirically unstable [53]. Beyond these empirical studies, one work rigorously defines CLIP as a mapping from images and texts to a latent unit hypersphere [6]. While Brody *et al.* focus on CLIP's inability to capture boolean logic, our work extends this perspective to the entire CLIP model and three complex reasoning tasks. Building on our analytical insights, we propose a simple solution to compute a more comprehensive similarity score.

Extended related work is provided in the Supplements.

## 3. The Ideal CLIP

The goal of CLIP is "to efficiently learn visual concepts from natural language supervision" [46]. What are *visual concepts*? We first establish formal notations with which to discuss visual concepts in the latent space. Once we have the tools to discuss CLIP mathematically, we will formalize the properties necessary to make CLIP ideal.

Fig. 2 illustrates how we think about the relationship between the world, its representation in images and language, and the eventual projection of images and text onto CLIP space. In this view, the "gist" of scene perception is broken down into categorization of objects, their salient attributes, and the spatial layout of the scene, as proposed by the cognitive psychology literature [5, 13, 23, 39, 49, 61]. Theories of perceptual symbol systems [3, 4, 17] propose that object categories and spatial relationships are atomic sym-
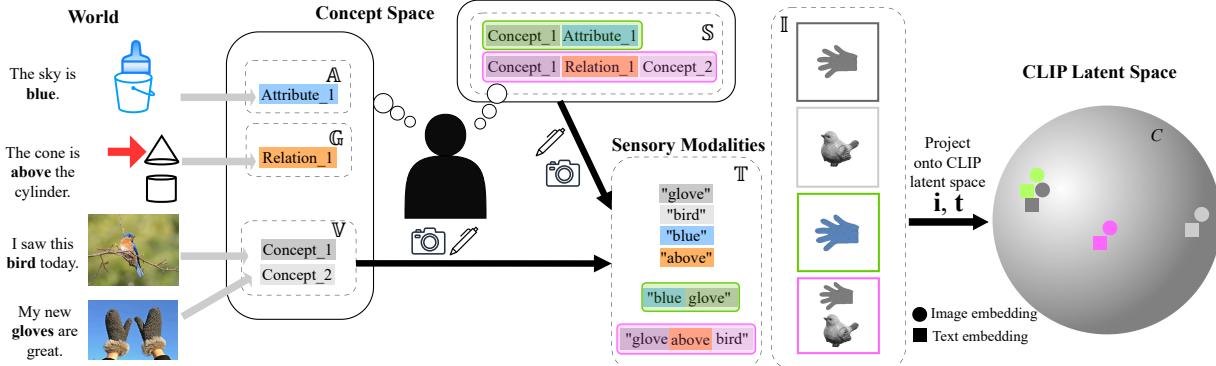
Figure 2. Graphical illustration of defined concept sets. Humans can parse visual stimuli from the real world and organize them into Object Concepts $\mathbb{V}$, Attributes which adorn objects $\mathbb{A}$, and Relationships between objects $\mathbb{G}$. These concepts can be ordered into Composed Scenes $\mathbb{S}$. Here, $\mathbb{V}, \mathbb{A}, \mathbb{G},$ and $\mathbb{S}$ are strict subsets of the set of all real world concepts. We can communicate these composed concepts via language $\mathbb{T}$, or by taking pictures of exemplars $\mathbb{I}$. These composed language or image modalities can be projected onto the CLIP latent space $C$ via a text encoder $\mathbf{t}$ or an image encoder $\mathbf{i}$. Elements in distinct sets with the same color have a one-to-one correspondence.

bols, combinatorially and recursively combined to create scene representations. The Dual Coding theory of cognition [41–43] proposes that these composed concepts can be represented as language, imagery, or both.

In the following we formalize these intuitions with precise definitions, with the goal of analyzing CLIP properties.

### 3.1. Concepts, attributes, and compositions

This section sets up a simplified world, which is a strict subset of the real world, within which we will explore whether it is possible to have an *ideal* CLIP. Our goal is to show that this is not possible and therefore, *a fortiori*, an ideal CLIP is not possible in more general worlds.

**Definition 1. (CLIP Embedding Space $C$).** CLIP embeddings are unit vectors which occupy an $N$-dimensional unit hypersphere. Let $\mathbb{I}$ denote any set of images and let $\mathbb{T}$ any set of texts. Then a CLIP-like model is a pair of functions $(\mathbf{i}, \mathbf{t})$ where for some Euclidean vector space $C$, consisting of $N$-dimensional unit vectors, we have that $\mathbf{i} : \mathbb{I} \to C$ and $\mathbf{t} : \mathbb{T} \to C$. That is, $\mathbf{i}$ and $\mathbf{t}$ are injective functions which map images and text descriptions into $C \in S^{N-1}$ where $S^{N-1}$ denotes the unit-($N$-1) sphere: $S^{N-1} = \left\{ \mathbf{i}(x), \mathbf{t}(x) \in \mathbb{R}^N : \|\mathbf{i}(x)\|_2 = \|\mathbf{t}(x)\|_2 = 1 \right\}$. Necessary properties will be detailed in the next section.

**Definition 2. (Atomic Concept Set $\mathbb{V}$).** In the real world, there exist $\gg N$ object concepts. These have hierarchy, and can be homonyms, polysemous, heteronyms, or synonyms. Let $\mathbb{V}$ be a small subset of this real world object concept set, where every $x$ in $\mathbb{V}$ has a distinct visual and linguistic counterpart, and all $x \in \mathbb{V}$ are semantically mutually exclusive, *e.g.* no $x$ is a subclass of another. We select $\mathbb{V}$ so that it contains fewer than $N$ objects: $|\mathbb{V}| = M \leq N$.

Each concept, *e.g.* "apple", corresponds to many physical objects, and each object may be portrayed in many different images. Now consider one representative element,

one image, for each set of images corresponding to a concept, and one representative element for each set of texts corresponding to a concept. In this way, we are defining injective functions $f_{V,I}: \mathbb{V} \to \mathbb{I}$ and $f_{V,T}: \mathbb{V} \to \mathbb{T}$ which act upon every element in $\mathbb{V}$. For the representative image which visually represents concept $x \in \mathbb{V}$, let its CLIP embedding be: $\mathbf{i}(f_{V,I}(x))$, which unless otherwise specified we will denote as $\mathbf{i}(x) \in S^{N-1}$. Similarly we denote some unique text representation of $x$ as $\mathbf{t}(f_{V,T}(x)) = \mathbf{t}(x) \in S^{N-1}$. Justifications for assumptions are provided in the supplements. Object concepts are shown in gray in Fig. 2.

**Definition 3. (Attribute Representation Set $\mathbb{A}$).** In the real world, there exist $\gg N$ attribute concepts. An attribute is not itself an object, but rather a property of the object, like color, size, and material. Then let $\mathbb{A}$ be a subset of this large attribute set, where every $a$ in $\mathbb{A}$ is a discrete and visually distinct attribute, and can be applied to every instance of $x$ in $\mathbb{V}$. For every $a \in \mathbb{A}$, $a$ has one equivalent item in $\mathbb{T}$, such that there exists some injective function $f_{A,T}: \mathbb{A} \to \mathbb{T}$.

Let $\mathbf{i}(x_a)$ denote an image embedding where concept $x$ has attribute $a$. For any given object, some attributes are highly characteristic (*e.g.,* "red" is a defining feature for an "apple"), some are relatively neutral (*e.g.,* "red" for a "car" may not be as informative), and some are unlikely to be associated with that object (*e.g.,* "red" for a "raccoon"). The representative image embedding $\mathbf{i}(x) \ \forall x \in \mathbb{V}$ implicitly portray the attributes that are most likely or typical for that concept. Attributes are shown in blue in Fig. 2.

**Definition 4. (Compositional Concept Set $\mathbb{G}$)** In the real world, there exist $\gg N$ compositional concepts, *i.e.,* concepts which describe the composition of a scene. Then let $\mathbb{G}$ be a subset of the real compositional concept set, consisting only of elements which describe the absolute location of one object (denoted $g^{<loc>}$) or the spatial relationship

3

| Notation | Explanation |
|---|---|
| $\mathbf{i}(\cdot), \mathbf{t}(\cdot)$ | CLIP image and text embeddings |
| $\mathbb{I}, \mathbb{T}$ | Set of images and texts |
| $\mathbb{V}$ | Set of atomic object concepts |
| $\mathbb{A}$ | Set of atomic attributes |
| $\mathbb{G}$ | Set of compositional relations |
| $\mathbb{S}$ | Finite ordered combinations of elts from $\mathbb{V}, \mathbb{A}, \mathbb{R}$ |
| $\mathbf{i}(x_a)$ | Image embedding of object $x$ with attribute $a$ |

Table 1. Summary of notations.

between two objects (denoted $g^{<rel>}$). Each element in $\mathbb{G}$ has a unique equivalent text description in $\mathbb{T}$. That is, there exists an injective function $f_{G,T} \colon \mathbb{G} \to \mathbb{T}$.

For example, for the compositional concept of one object being above another object, we use the notation $g_{above}^{<rel>}$, which has the text representation "above". (In the next definition we explain our notation for using compositional concepts in context with other concepts.) Compositional concepts are shown in orange in Fig. 2. When two objects $x$ and $y$ appear in the same image, a spatial relationship is there between them. If one wishes to marginalize this aspect out one may use the notation $\mathbf{i}(x, y)$ to denote the mean embedding of images depicting concepts $x$ and $y$, averaged over all spatial arrangements of $x$ and $y$ in the fixed frame.

**Definition 5. (Combinatorial Clause Set $\mathbb{S}$)** In the real world, object, attribute, and compositional concepts can be organized to create semantic *clauses*. These are a combination of concepts which add context to a unified scene. Let $\mathbb{S}^*$ be a set of all finite ordered combinations of elements from $\mathbb{V}$, $\mathbb{A}$, and $\mathbb{G}$. It is a union of all $n$-fold Cartesian products of the full outer joined set of the three.

$$\mathbb{S}^* = \bigcup_{n=1}^{\infty} \left( \mathbb{A} \cup \mathbb{G} \cup \mathbb{V} \right)^n$$

Every element in $\mathbb{S}^*$ has an equivalent text description in $\mathbb{T}$. Now, let $\mathbb{S}$ be a subset of elements from $\mathbb{S}^*$ whose textual counterpart is grammatically correct, such that there exist equivalent and distinct mappings between every item in $\mathbb{S}$ and $\mathbb{I}, \mathbb{T}$. This means there exist injective functions $f_{S,I} \colon \mathbb{S} \to \mathbb{I}$ and $f_{S,T} \colon \mathbb{S} \to \mathbb{T}$. For example, the joined clause of the object concept of a glove and the attribute concept of the color blue, $s = (x_{glove}, a_{blue}) \in \mathbb{S}$, maps to the textual phrase "blue glove" and to an image of a blue glove.

By convention, any $g^{<loc>}$ declares the location of the preceding object in the clause. That is, $f_{S,T}((x, g_{right})) = $ "$x$ to the right". Similarly, any $g^{<rel>}$ in between two objects influences both in that order: $f_{S,T}((x, g_{below}, y)) = $ "$x$ below $y$". Combinations of concepts are shown in green in Fig. 2. For simplicity, we denote $\mathbf{i}(f_{S,I}((x, a))) = \mathbf{i}(x_a)$ unless otherwise specified. Similarly, we denote scenes composed of multiple objects as $\mathbf{i}(f_{S,I}((x, y))) = \mathbf{i}(x, y), \forall x, y \in \mathbb{V}$ as shorthand.

## 3.2. Geometric Requirements for an Ideal CLIP

We will assume that CLIP encoders can project any image and text as vectors in an "ideal" location in the latent space *iff* there exists such a location. By "ideal location", we mean the semantic distance between image and text embeddings aligns well with human understanding. If under this oracle encoder setting we find inherent contradictions which prohibit an "ideal" latent space from existing, the conclusions will generalize to the real life setting which is a superset of the concepts we consider, with non-oracle $\mathbf{i}, \mathbf{t}$. We justify the need for each condition in the supplements. $C$ is "ideal" if it fulfills all following conditions:

**Condition 1. (Concept Categorization)** Satisfaction of this condition requires that (1.1) $C$ represents basic descriptions and image content.

$$\mathbf{i}(x) \cdot \mathbf{t}(x) \;>\; \mathbf{i}(x) \cdot \mathbf{t}(y)$$
$$\mathbf{i}(x, y) \cdot \mathbf{t}(x) \;>\; \mathbf{i}(x, y) \cdot \mathbf{t}(z) \quad \forall \, x, y, z \in \mathbb{V}$$

(1.2) Images that contain the same semantic concept(s) but differ due to an attribute or scene composition, should have higher cosine similarity with each other than with an image that contains a different set of semantic concepts.

$$\mathbf{i}(x_a) \cdot \mathbf{i}(x_b) \;>\; \mathbf{i}(x_a) \cdot \mathbf{i}(y)$$
$$\mathbf{i}(x, g_1^{<loc>}) \cdot \mathbf{i}(x, g_2^{<loc>}) \;>\; \mathbf{i}(x) \cdot \mathbf{i}(y)$$
$$\forall \, x, y \in \mathbb{V}, \forall \, a, b \in \mathbb{A}, \qquad \forall \, g_1^{<loc>}, g_2^{<loc>} \in \mathbb{G}$$

**Condition 2. (Attribute Binding)** $C$ respects attribute binding. More specifically: (2.1) concepts with different attributes are not parallel in CLIP space.

$$\mathbf{i}(x_a) \cdot \mathbf{i}(x_b) < 1 \quad \forall a, b \in \mathbb{A}$$

(2.2) Images representing a concept with a specific attribute are closer in CLIP space to its text embedding.

$$\mathbf{i}(x_a) \cdot \mathbf{t}(a) > \mathbf{i}(x_b) \cdot \mathbf{t}(a)$$

(2.3) Images with the same concepts and attributes present but in different pairings are not parallel in CLIP space.

$$\mathbf{i}(x_a, y_b) \cdot \mathbf{i}(x_b, y_a) < 1$$

**Condition 3. (Spatial Relationship)** C respects spatial locations or relationships of objects. This requires that (3.1) images where the same object is in a different location must not have identical embeddings.

$$\mathbf{i}(x, g_1^{<loc>}) \cdot \mathbf{i}(x, g_2^{<loc>}) < 1, \quad \forall g_1^{<loc>}, g_2^{<loc>} \in \mathbb{G}$$

(3.2) Images with the same objects but in different spatial relationships must not have identical embeddings.

$$\mathbf{i}(x, g_3^{<rel>}, y) \cdot \mathbf{i}(x, g_4^{<rel>}, y) < 1, \quad \forall g_3^{<rel>}, g_4^{<rel>} \in \mathbb{G}$$

(3.3) Images where an object is in the same location or relationship must be semantically closer than images where it is in a different location or relationship.

$$\mathbf{i}(x, g_1, y) \cdot \mathbf{i}(x, g_1, z) > \mathbf{i}(x, g_1, y) \cdot \mathbf{i}(x, g_2, z)$$

**Condition 4. (Negation)** $C$ respects negation. This requires that (4.1) text embeddings and their negated counterparts must have a similarity score lower than any other

pairs.
$$\mathbf{t}(x) \cdot \mathbf{t}(\neg x) < \mathbf{t}(y) \cdot \mathbf{t}(\neg x), \quad \forall\, x, y \in \mathbb{T}$$

(4.2) The positive version for one concept must have a lower cosine similarity score with any other concept than the negated counterpart.
$$\mathbf{i}(x) \cdot \mathbf{t}(\neg x) < \mathbf{i}(x) \cdot \mathbf{t}(y)$$

(4.3) Two distinct negated concepts must have higher cosine similarity than distinct non-negated concepts, as they have greater semantic overlap.*
$$\mathbf{t}(x) \cdot \mathbf{t}(y) < \mathbf{t}(\neg y) \cdot \mathbf{t}(\neg x)$$

## 4. Geometric Contradictions in CLIP

We aim to prove that if $C$ meets Condition 1, it cannot meet any other condition to be ideal. Therefore no $C$ exists which is ideal. Due to space limitations, we only illustrate in detail how fulfillment of Condition 1 prohibits fulfillment of Condition 2 then summarize the rest. Please refer to the Supplements for detailed proofs of the other conditions.

### 4.1. Contradiction for Conditions 1 and 2

**Lemma 1. Embeddings of images or texts with two obect concepts must be a linear superposition of the respective single object concept embeddings.**

*Proof.* We derive the ideal placement for $\mathbf{i}(x^1, x^2)$ to satisfy Condition 1.1. For this proof, instead of $x, y$ we denote distinct concepts with superscripts: $\mathbf{i}(x^1), \mathbf{i}(x^2), \ldots$ for $k \in [1, M]$ (to avoid confusion with attributes, which are denoted as subscripts.) The condition states $\mathbf{i}(x^1, x^2)$ must have high cosine similarity with $\mathbf{i}(x^1)$ and $\mathbf{i}(x^2)$, and low cosine similarity with all other $\mathbf{i}(x^j)$. More formally, $C$ must solve the following optimization problem:

$$\mathbf{i}(x^1, x^2) = \operatorname*{argmax}_{\mathbf{i}(x^1, x^2)} \Big[ \mathbf{i}(x^1, x^2) \cdot \mathbf{i}(x^1) + \mathbf{i}(x^1, x^2) \cdot \mathbf{i}(x^2)$$
$$- \sum_{j=3}^{M} \mathbf{i}(x^1, x^2) \cdot \mathbf{i}(x^j) \Big] \ \text{ s.t. } \ \|\mathbf{i}(x^1, x^2)\| = 1 \quad (1)$$

Here, the first two terms guide the local placement of $\mathbf{i}(x^1, x^2)$, while the last term introduces a global constraint to avoid proximity to other embeddings. The constraint ensures that all embeddings must lie on the unit hypersphere. We can expand the sum to see that:

$$\mathbf{i}(x^1, x^2) = \operatorname*{argmax}_{\mathbf{i}(x^1, x^2)} \Big[ \mathbf{i}(x^1, x^2)$$
$$\cdot \Big( \mathbf{i}(x^1) + \mathbf{i}(x^2) + \mathbf{i}(x^1) + \mathbf{i}(x^2) - \sum_{j=1}^{M} \mathbf{i}(x^j) \Big) \Big] \quad (2)$$

Since random vectors in high dimensions will be approximately symmetrically distributed, $\sum_{j=1}^{M} \mathbf{i}(x^j) \approx 0$. The optimum is then reached when $\mathbf{i}(x^1, x^2)$ is parallel to $\mathbf{i}(x^1) +$

---

*If $|\mathbb{V}| = M$, "not $x$" and "not $y$" correctly describe $M - 2$ concepts, while "$x$" and "$y$" describe none of the same things.

$\mathbf{i}(x^2)$. Thus we see $\mathbf{i}(x^1, x^2)$ is a normalized superposition of $\mathbf{i}(x^1)$ and $\mathbf{i}(x^2)$, and lies on the geodesic arc between $\mathbf{i}(x^1)$ and $\mathbf{i}(x^2)$ on the hypersphere, i.e.,

$$\mathbf{i}(x^1, x^2) = \frac{\mathbf{i}(x^1) + \mathbf{i}(x^2)}{\|\mathbf{i}(x^1) + \mathbf{i}(x^2)\|} \quad (3)$$

$\square$

**Lemma 2. $C$ cannot distinguish between different attribute bindings. That is, $\mathbf{i}(x_a, y_b) = \mathbf{i}(x_b, y_a)$.**

*Proof.* We start by deriving the ideal location for $\mathbf{i}(x_a)$ in $C$ where the distance between this unit vector and $\mathbf{i}(x)$ and $\mathbf{t}(a)$ is strictly semantically correct per Conditions 1 and 2.1-2.2. Then we derive $\mathbf{i}(x_a, y_b)$ and $\mathbf{i}(x_b, y_a)$ to meet Condition 1.1. We will find that this derivation contradicts Condition 2.3.

Condition 1.2 states that images of the same object that differ only in attributes (*e.g.,* red car, black car) should be more similar to each other in $C$ than than images of distinct objects. This means that the attribute-specific image embedding can be expressed as a small perturbation of the representative image embedding for that object:

$$\mathbf{i}(x_a) = (1 - \delta)\mathbf{i}(x) + \mathbf{v} \quad (4)$$

where $\delta \ll 1$ is some small positive constant and vector $\mathbf{v}$ denotes the attribute-specific change such that $\|\mathbf{i}(x_a)\| = 1$ and $\mathbf{i}(x_a) \cdot \mathbf{i}(x) \geq 1 - \delta$.

Now we consider two objects $x$ and $y$ and two attributes $a$ and $b$ where the objects are agnostic to the attributes. More precisely, we consider $x, y, a, b$ where the semantic distance between each attribute's text embedding and both objects is equal, *i.e.,* $\mathbf{i}(x) \cdot \mathbf{t}(a) = \mathbf{i}(y) \cdot \mathbf{t}(a) = \cos(\theta)$ and $\mathbf{i}(x) \cdot \mathbf{t}(b) = \mathbf{i}(y) \cdot \mathbf{t}(b) = \cos(\omega)$. A concept quartet that satisfies this criterion could be $x$:car, $y$:ball, $a$:red, $b$: black.

Condition 2.2 states that $\mathbf{i}(x_a) \cdot \mathbf{t}(a) > \mathbf{i}(x) \cdot \mathbf{t}(a)$. Therefore, a large component of $\mathbf{v}$ in Eq. (4) must be in the direction of $\mathbf{t}(a)$. Below, we will show that if we make the strong assumption that $v = p\mathbf{t}(a)$, then $\mathbf{i}(x_a, y_b) = \mathbf{i}(x_b, y_a)$. Then we show that the result will be robust even if $\mathbf{v} = p\mathbf{t}(a) + \boldsymbol{\epsilon}$ for some noise vector $\boldsymbol{\epsilon}$. Letting $\mathbf{v} = p\mathbf{t}(a)$, we write $\mathbf{i}(x_a)$ as a superposition:

$$\mathbf{i}(x_a) = (1 - \delta)\mathbf{i}(x) + p\mathbf{t}(a) \quad (5)$$

We derive $p$ in Sec E.4 of the Supplements as:
$$p = -(1 - \delta)\cos\theta \pm \tfrac{1}{2}\sqrt{4(1 - \delta)^2 \cos^2\theta + 8\delta - 4\delta^2} \quad (6)$$

Notice that $p$ only depends on $\theta$, where $\theta = \arccos(\mathbf{i}(x) \cdot \mathbf{t}(a)) = \arccos(\mathbf{i}(y) \cdot \mathbf{t}(a))$. By the same reasoning as above, for object $y$ with attribute $a$, we have:

$$\mathbf{i}(y_a) = (1 - \delta)\mathbf{i}(y) + p\mathbf{t}(a) \quad (7)$$

Similarly for attribute $b$, $\mathbf{i}(x_b)$ and $\mathbf{i}(y_b)$ share the same weighting factor $q$ for $\mathbf{t}(b)$ as follows:

$$\mathbf{i}(x_b) = (1 - \delta)\mathbf{i}(x) + q\mathbf{t}(b)$$
$$\mathbf{i}(y_b) = (1 - \delta)\mathbf{i}(y) + q\mathbf{t}(b) \quad (8)$$

5

Now consider the composite image embedding $\mathbf{i}(x_a, y_b)$. Per Eq. (3), it is decomposed into a superposition:

$$\mathbf{i}(x_a, y_b) = \frac{(1-\delta)(\mathbf{i}(x) + \mathbf{i}(y)) + p\mathbf{t}(a) + q\mathbf{t}(b)}{2} \qquad (9)$$
$$= \mathbf{i}(x_b, y_a)$$

This shows that the composite embedding is identical regardless of which object is paired with which attribute. In other words, the distinct attribute bindings become indistinguishable in the final embedding. This equivalence violates Condition 2.3, which requires that different attribute-object bindings produce distinct embeddings.

For cases where $\mathbf{v} = p\mathbf{t}(a) + \boldsymbol{\epsilon}$ for some random perturbation vector $\boldsymbol{\epsilon}$, we show in Sec E.3. of the Supplements that the results in Eq. (9) are robust to the addition of $\boldsymbol{\epsilon}$. $\square$

### 4.2. Contradictions to Other Conditions

Now we proceed to illustrate how the linear superposition implied by Condition 1 contradicts other conditions. If Condition 1 is fulfilled, then the following impossibilities occur:

- **Condition 3:** Spatial relationships between objects must embed as $\mathbf{i}(x, g^{<rel>}, y) = (1-\delta)\,\mathbf{i}(x, y) + \mathbf{e}_\perp$ per Condition 1.2. We show that for a simple case of 3 images, where the same localization is present in image 1 and 2 and the same relationship is present in image 1 and 3, we encounter $\mathbf{e}_{\perp,1} \cdot \mathbf{e}_{\perp,2} = -\mathbf{e}_{\perp,1} \cdot \mathbf{e}_{\perp,3}$. This means Conditions 3.1 and 3.2 cannot be simultaneously satisfied.
- **Condition 4:** We find that satisfaction of Conditions 4.1 and 2 in $C$ necessitates $\mathbf{t}(\neg x) = -\mathbf{t}(x)$. This produces the following effect: $\mathbf{t}(\neg x^j) \cdot \mathbf{t}(x^k) > \mathbf{t}(\neg x^j) \cdot \mathbf{t}(\neg x^k)$, violating Condition 4.3.

Please refer to the Supplements for detailed proofs.

## 5. Rescuing the CLIP Latent Space

Is CLIP beyond rescue, or can its learned embeddings be improved? Although its latent space lacks compositional expressivity, its ease of use and powerful image-text organization are undeniable. Hence, we investigate methods to make use of the rich information learned by CLIP for a more principled evaluation of the text to image semantic distance.

Our findings indicate that CLIP is structurally flawed. Re-training, fine-tuning, or simply re-projecting CLIP embeddings will still yield a latent space that lacks the desired properties. Similarly, applying an alternative analytical or learned scoring mechanism on top of the learned embeddings cannot work if two distinct images and texts embed to the same location in $C$ per Lemma 2.

***Thus: any fix must alter the fact that CLIP represents images and texts as unit vectors.*** Rather than designing an approach that imposes re-training a full model, we explore an extension to CLIP that is based on the existing CLIP encoders. More specifically: we explore a solution based on three ideas. First, retain the full token and image patch embeddings from CLIP. Second, score matches through a learned mechanism, rather than using cosine similarity. Third, introduce constant, rather than learned, representations of spatial relationship words. These ideas are discussed in detail in the next section.

### 5.1. Dense Cosine Similarity Maps

Instead of deriving a single-point CLIP embedding from EOS token from the text embeddings and the CLS token from the image embeddings, we propose to compute the pairwise cosine similarity between all text tokens and all image patches, and then to use a learned scoring mechanism on the resulting dense cosine similarity map (DCSM). The intuition is that transformers [57], in learning the correct CLS and EOS token pair, store useful information in the token and patch level embeddings. While dense image patches have been used in tasks like image segmentation or dense label generation [48, 73], our approach is the first to extract embeddings at both the token and patch level, and densely compute the cosine similarity across the two representations to formulate the task of image-text pair scoring as a pattern recognition problem.

Fig. 3 depicts dense DCSMs for four distinct sentences and one image prompt. The columns have been rearranged to cluster background (with CLS patch being the first column), backpack, and glove patches. Our DCSMs clearly encode object and attribute locality. As expected from the contrastive training procedure, the highest score in each DCSM is between the EOS token and CLS patch embedding. Interestingly, text tokens describing a specific object have low cosine similarity scores with the patches containing that object. We hypothesize that this occurs because background patches store global information about the im-
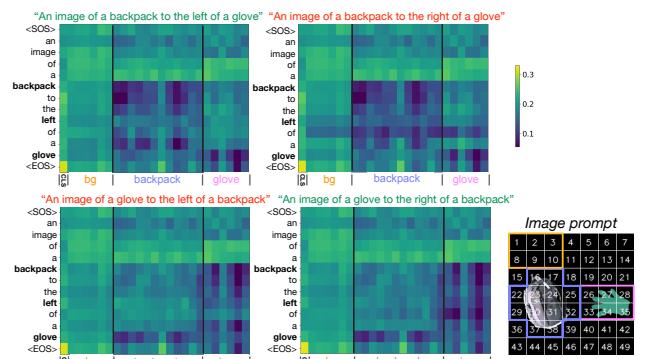


Figure 3. **Empirical Dense Cosine Similarity Maps.** Each one of four matrices shows a DCSM between a different sentence and the same pictured image. For each subfigure: the y axis shows the text tokens and the x axis varies by image patch. We cluster the patches by region as shown in the image. Each pixel value is the cosine similarity score between that token and patch embedding. Green sentences correctly represent the image and red ones do not.
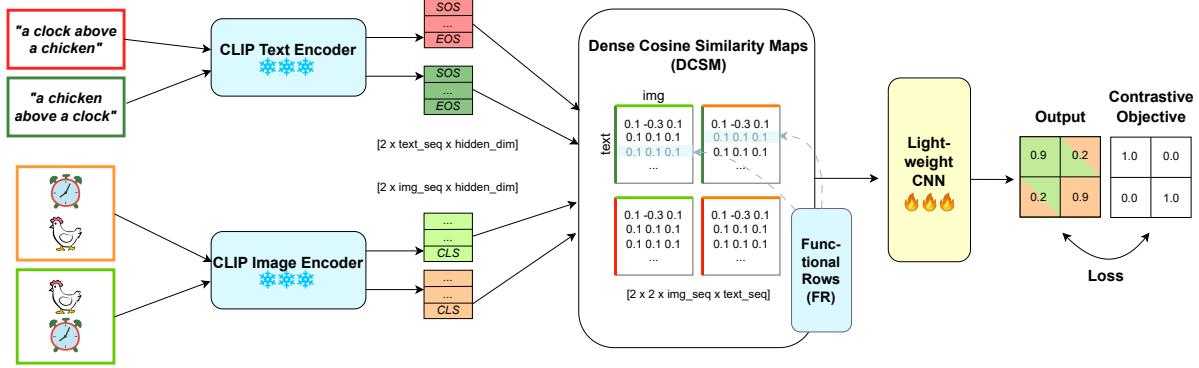
Figure 4. **Schematics of our proposed pipeline and training of its scoring function.** Every sample seen during training contains one hard positive caption and image pair, and a hard negative caption and image pair. Images and texts are passed through frozen CLIP encoders to compute the DCSMs, and then functional rows (FRs) for compositional words are inserted before the DCSMs are scored.

age, including information about what objects are present, while the salient patches themselves preserve local information, making them more distant from text embeddings in general. This observation is analogous to how vision transformers (VIT) use non-salient patches as "register tokens" [11] to store higher-level global information.

**Functional Rows (FR).** We observe that for both correct sentence-image pairs, the DCSMs often look very similar. Much like background patches in an image, text embeddings for compositional concepts, which we will refer to as *functional words* (*e.g.,* "left of","right of") which lack immediate visual counterparts only contribute spurious information to the DCSM as the joint encoder training emphasizes words with direct visual references.

This motivates our proposed *functional rows* for DCSMs. Specifically, we discard the original rows containing functional words in DCSMs, and replace them with a constant vector with randomly chosen entries.

**Model Pipeline** (see Fig. 4)**.** Unlike classical CLIP scoring which discards all but the EOS and CLS tokens to produce a

single scalar, our proposed method retains all token embeddings and projects them into a joint latent space to compute the DCSM. For all experiments, we use a lightweight CNN with 2 convolutional layers and a hidden dimension of 128, yielding a 20-fold reduction in parameters.

**Training.** The network is trained with a batch size of 8, which is 4000 times smaller than CLIP's original 32,768. Our total training data per model is around 20,000 samples, which is actually 1.5 times smaller than the mini-batch size of CLIP. For a small set of functional words (elements of $f_{G,T}(g)$ for all $g \in \mathbb{G}$), we keep a dictionary of unique fixed FRs and overwrite the corresponding rows in the DCSM. As illustrated in Fig. 5, in an idealized setting, DCSMs with FRs allow for clear disambiguation between correct and incorrect image-text pairs, reducing the task of parsing visuolinguistic semantics to pattern recognition. Notably, our model does *not* see any input image or text at all; it simply learns to recognize the syntactic patterns from the DCSMs.

We train one model on synthetic data created from Objaverse [12] renders and another on a subset of COCO2017. Additional training details are provided in the Supplements.

### 5.2. Performance Evaluation

We compare our method against a range of models which utilize a joint vision-language embedding space. Namely: CLIP from OpenAI [46] and OpenCLIP [9], NegCLIP [66] which is a finetuned version of OpenCLIP with mined hard negatives from all of COCO-train, Coca [65] which is an image-text encoder-decoder model trained with contrastive and captioning loss , and BLIP [25] which bootstraps noisy web scraped data with synthetic, filtered captions. The models are evaluated on datasets of three categories: CLEVR-bind [22], a composite form of Natural Colors Dataset (NCD) [2] per [63], and VG_attribution from Attribute, Relation, Order (ARO) [66] for attribute binding; WhatsUp, COCO-QA, and VG-QA [18] for spatial relationships and localization; and NegBench [1] for negation. Table 2 shows the results.
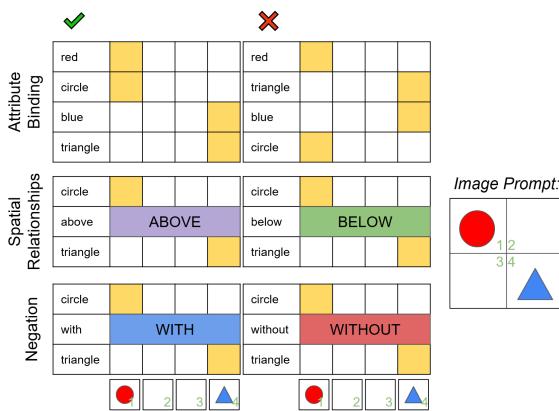


Figure 5. Graphic illustration of retained topology in DCSMs. The image prompt has 4 patches, and each DCSM shows the dense cosine similarity between each of its patches and the text tokens of the text prompts. Sentences under the green checkmark are a correct pair with the image, while those with a red X are incorrect.

| | Attribute Binding | | | Spatial Reasoning | | | Negation | |
|---|---|---|---|---|---|---|---|---|
| **Model Name** | $\mathbf{CLEVR}_{bind}$ | **NCD** | **VG_attr** | **WhatsUp** | $\mathbf{COCO}_{1\&2obj}$ | $\mathbf{VG}_{1\&2obj}$ | $\mathbf{NegBench}_{coco}$ | $\mathbf{NB}_{voc}$ |
| $CLIP_{ViTB/32}$ | 22.2 | 71.3 | 61.3 | 31.9 | 47.0 | 47.1 | 39.2 | 38.3 |
| $CLIP_{ViTB/16}$ | 20.2 | 63.1 | 61.8 | 30.5 | 48.9 | 51.5 | 41.5 | 37.9 |
| NegCLIP | 21.8 | 79.2 | **72.4** | 33.2 | 46.1 | 47.2 | 31.4 | 26.5 |
| CoCa | 19.2 | 48.7 | 50.8 | 24.5 | 48.6 | 49.5 | 21.6 | 20.8 |
| BLIP | 11.1 | 56.2 | 59.4 | 24.4 | 48.5 | 50.4 | 20.5 | 20.7 |
| SigLIP | 13.3 | 53.6 | 48.4 | 26.0 | 47.4 | 51.1 | 26.3 | 29.7 |
| $\mathbf{DCSM}_{synth}$ | 31.0 | 93.6 | 60.9 | 62.6 | 65.6 | 64.4 | 38.8 | 35.2 |
| $\mathbf{DCSM}_{coco}$ | **39.9** | **95.7** | 68.1 | **63.7** | **72.4** | **67.0** | **48.6** | **49.0** |
| Random Chance | 25.0 | 50.0 | 50.0 | 25.0 | 50.0 | 50.0 | 25.0 | 25.0 |

Table 2. Performance comparison across different models on various benchmarks. The top row categorizes each benchmark into the condition being addressed. Our DCSM takes CLIP ViTB/16 as the base model. Best scores for each dataset are bolded. Evaluation datasets contain scene compositions and attribute, spatial, and negation words not included in the training set for out model.

Our method significantly outperforms CLIP-like baselines, likely due to the increased dimensionality of DCSMs, which replace scalar cosine similarity with a dense topological map. Specifically, patch indices on the image embedding axis store spatial information, while the token indices on the text embedding axis preserve the semantic ordering. DCSMs trained on COCO outperform the baselines across all datasets except on VG-attribution [66], where NegCLIP benefits from similar hard negatives. Improvements on NegBench are lower, likely due to the data gap between our templated training captions and the natural language captions of this evaluation set.

Notably, our model generalizes well to unseen attribute, spatial, and negation concepts despite the limited training set. In particular, our model is only trained on templated two-object captions but works well for single object captions on CLEVR-bind and subsets of COCO and VG-spatial. This suggests the downstream network learns syntactic patterns rather than overfitting to the training templates. Beyond performance gains, DCSMs offer greater interpretability. Unlike naive CLIP embeddings whose actual values can be arbitrary, DCSMs are human interpretable (see Fig. 3), which makes downstream usage more intuitive.

For detailed ablation studies on the impact of dense maps, FRs, classification performance, and scaling analysis, please refer to Sec. C of the Supplements.

### 5.3. Generalization to Natural Language

To test the generalizability of DCSMs to open vocabulary settings, we incorporate LLMs-in-the-loop to dynamically update the lookup table for FRs and reformat natural language sentences for more compact DCSMs. We explore this in a toy setting by training our lightweight CNN from scratch with either 5k or 10k COCO images whose captions have had their nouns swapped. We choose this particular type of intervention as CLIP-like VLM performance across the board was lowest for this category of hard negatives in Sugarcrepe [15]. On top of the training pipeline described

| Model | $VLC_{vg\ spatial}$ | $SC_{swap\ obj}$ |
|---|---|---|
| $CLIP_{ViTB/16}$ | 50.8 | 60.2 |
| $CLIP_{ft.synth}$ | 56.9 | 50.0 |
| $\mathbf{DCSM}_{synth}$ | 50.6 | 59.8 |
| $\mathbf{DCSM}_{open\ vocab}$ | **63.5** | **63.8** |

Table 3. Sugarcrepe results with LLM prompt simplification are averaged across 6 trials. VL-Checklist captions were not simplified at test time. We report the higher score of the two trained models for the open-vocabulary DCSM.

in Sec. 5.1, we prompt gpt-4o-mini to extract newly encountered functional words in the training data. We evaluate this paradigm on the swap-object split of Sugarcrepe and the VG-spatial split of VL-Checklist [72], against naive and finetuned CLIP. Details for training and prompting are in the Supplements. Results are shown in Table 3.

VL-checklist captions are short but contain novel functional words not present in the static lookup table. Sugarcrepe captions are longer, with complex syntactic arrangements as well as novel functional words. We see that the finetuned CLIP model overfits to simple clauses and fails at the more complex setting, while the closed vocabulary DCSM performance is consistent with that of naive CLIP. The open vocabulary DCSM shows a modest performance increase in both settings. This suggests that LLM for dynamic FR update is a promising avenue, though long winded prompts are still a bottleneck for the lightweight CNN model. We predict that increasing the dataset size and quality, as well as improving scoring model complexity will help refine natural language performance.

## 6. Conclusions

Inherent geometric limitations in CLIP prevent the correct representation of image content, attribute binding, spatial localization, and negation. To address this problem we tap the information from CLIP's text and image encoders, creating a richer latent space with Dense Cosine Similarity Maps and Functional Rows. We evaluate our simple so-

lution against a wide array of benchmarks and find that it surpasses SoTA performance. Future work includes scaling up our prototyped pipeline to expand on its open vocabulary potential with more sophisticated CLIP base models, scoring models, LLMs, and training data. Further analyses to discover other fundamental restrictions in existing VLM architectures will help guide future design choices.

# References

[1] Kumail Alhamoud, Shaden Alshammari, Yonglong Tian, Guohao Li, Philip Torr, Yoon Kim, and Marzyeh Ghassemi. Vision-language models do not understand negation, 2025. 1, 2, 7, 12

[2] Saeed Anwar, Muhammad Tahir, Chongyi Li, Ajmal Mian, Fahad Shahbaz Khan, and Abdul Wahab Muzaffar. Image Colorization: A Survey and Dataset, 2024. arXiv:2008.10774 [cs]. 7

[3] Lawrence W Barsalou. Perceptual symbol systems. *Behav. Brain Sci.*, 22(4):577–660, 1999. 2

[4] Lawrence W Barsalou. Abstraction in perceptual symbol systems. *Philos. Trans. R. Soc. Lond. B Biol. Sci.*, 358(1435): 1177–1187, 2003. 2

[5] Irving Biederman. Recognition-by-components: a theory of human image understanding. *Psychol. Rev.*, 94(2):115–147, 1987. 2

[6] Justin Brody. On the Potential of CLIP for Compositional Logical Reasoning, 2023. 2

[7] Declan Campbell, Sunayana Rane, Tyler Giallanza, Nicolò De Sabbata, Kia Ghods, Amogh Joshi, Alexander Ku, Steven M. Frankland, Thomas L. Griffiths, Jonathan D. Cohen, and Taylor W. Webb. Understanding the Limits of Vision Language Models Through the Lens of the Binding Problem, 2024. arXiv:2411.00238 [cs]. 1, 2, 12

[8] Boyuan Chen, Zhuo Xu, Sean Kirmani, Brian Ichter, Danny Driess, Pete Florence, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. SpatialVLM: Endowing Vision-Language Models with Spatial Reasoning Capabilities, 2024. arXiv:2401.12168 [cs]. 2, 12

[9] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2818–2829, 2023. 7

[10] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023. 15

[11] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision Transformers Need Registers, 2023. arXiv:2309.16588 [cs]. 7

[12] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A Universe of Annotated 3D Objects, 2022. arXiv:2212.08051 [cs]. 7, 13

[13] Li Fei-Fei, Asha Iyer, Christof Koch, and Pietro Perona. What do we perceive in a glance of a real-world scene? *Journal of vision*, 7(1):10–10, 2007. 2

[14] Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning without training, 2022. 2, 12

[15] Cheng-Yu Hsieh, Jieyu Zhang, Zixian Ma, Aniruddha Kembhavi, and Ranjay Krishna. SugarCrepe: Fixing Hackable Benchmarks for Vision-Language Compositionality, 2023. arXiv:2306.14610 [cs]. 2, 8, 12, 19, 21

[16] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning*, pages 4904–4916. PMLR, 2021. 2, 12

[17] D. Joyce, L. Richards, A. Cangelosi, and K.R. Coventry. *On the foundations of perceptual symbol systems:: Specifying embodied representations via connectionism*, pages 147–152. Universitätsverlag Bamberg, Germany, 2003. 2

[18] Amita Kamath, Jack Hessel, and Kai-Wei Chang. What's "up" with vision-language models? Investigating their struggle with spatial reasoning, 2023. arXiv:2310.19785 [cs]. 1, 2, 7

[19] Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. DiffusionCLIP: Text-Guided Diffusion Models for Robust Image Manipulation, 2022. arXiv:2110.02711 [cs]. 1

[20] Soroush Abbasi Koohpayegani, Anuj Singh, K. L. Navaneet, Hadi Jamali-Rad, and Hamed Pirsiavash. GeNIe: Generative Hard Negative Images Through Diffusion, 2024. arXiv:2312.02548 [cs]. 1, 12

[21] Meir Yossef Levi and Guy Gilboa. The Double-Ellipsoid Geometry of CLIP, 2024. arXiv:2411.14517. 2

[22] Martha Lewis, Nihal V. Nayak, Peilin Yu, Qinan Yu, Jack Merullo, Stephen H. Bach, and Ellie Pavlick. Does CLIP Bind Concepts? Probing Compositionality in Large Image Models, 2024. arXiv:2212.10537 [cs]. 1, 2, 7, 12

[23] Fei Fei Li, Rufin VanRullen, Christof Koch, and Pietro Perona. Rapid natural scene categorization in the near absence of attention. *Proceedings of the National Academy of Sciences*, 99(14):9596–9601, 2002. 2

[24] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in neural information processing systems*, 34:9694–9705, 2021. 2, 12

[25] Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. BLIP: bootstrapping language-image pre-training for unified vision-language understanding and generation. *CoRR*, abs/2201.12086, 2022. 1, 7, 15

[26] Junyan Li, Delin Chen, Yining Hong, Zhenfang Chen, Peihao Chen, Yikang Shen, and Chuang Gan. CoVLM: Composing Visual Entities and Relationships in Large Language Models Via Communicative Decoding, 2023. arXiv:2311.03354 null. 2, 12

[27] Wei Li, Zhen Huang, Xinmei Tian, Le Lu, Houqiang Li, Xu Shen, and Jieping Ye. Interpretable Composition Attribu-

tion Enhancement for Visio-linguistic Compositional Understanding. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 14616–14632, Miami, Florida, USA, 2024. Association for Computational Linguistics. 2, 12

[28] Zejian Li, Chenye Meng, Yize Li, Ling Yang, Shengyuan Zhang, Jiarui Ma, Jiayi Li, Guang Yang, Changyuan Yang, Zhiyuan Yang, Jinxiong Chang, and Lingyun Sun. LAION-SG: An Enhanced Large-Scale Dataset for Training Complex Image-Text Models with Structural Annotations, 2024. arXiv:2412.08580 [cs] version: 1. 1, 12

[29] Weixin Liang, Yuhui Zhang, Yongchan Kwon, Serena Yeung, and James Zou. Mind the gap: Understanding the modality gap in multi-modal contrastive representation learning. 2

[30] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014. 13

[31] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024. 2, 12

[32] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 2, 12, 15

[33] Chengcheng Ma, Yang Liu, Jiankang Deng, Lingxi Xie, Weiming Dong, and Changsheng Xu. Understanding and Mitigating Overfitting in Prompt Tuning for Vision-Language Models, 2023. arXiv:2211.02219 [cs]. 12

[34] Damiano Marsili, Rohun Agrawal, Yisong Yue, and Georgia Gkioxari. Visual agentic ai for spatial reasoning with a dynamic api, 2025. 2, 12

[35] Sachit Menon and Carl Vondrick. Visual Classification via Description from Large Language Models. 2022. 2, 12

[36] Norman Mu, Alexander Kirillov, David Wagner, and Saining Xie. Slip: Self-supervision meets language-image pretraining. In *European conference on computer vision*, pages 529–544. Springer, 2022. 2, 12

[37] Kaleb Newman, Shijie Wang, Yuan Zang, David Heffren, and Chen Sun. Do Pre-trained Vision-Language Models Encode Object States?, 2024. arXiv:2409.10488 [cs] version: 1. 1, 2, 12

[38] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models, 2022. arXiv:2112.10741 [cs]. 1

[39] Aude Oliva. *Gist of the Scene*, page 251–256. Elsevier, 2005. 2

[40] Roni Paiss, Ariel Ephrat, Omer Tov, Shiran Zada, Inbar Mosseri, Michal Irani, and Tali Dekel. Teaching CLIP to Count to Ten, 2023. arXiv:2302.12066. 1, 12

[41] Allan Paivio. *Mental Representations*. Oxford University Press, 1990. 3

[42] A. Paivio. *Imagery and Verbal Processes*. Psychology Press, 2013.

[43] Allan Paivo and Ian Begg. *Psychology of language*. Prentice Hall, Old Tappan, NJ, 1981. 3

[44] Harold R Parks and Dean C Wills. An elementary calculation of the dihedral angle of the regular n-simplex. *Am. Math. Mon.*, 109(8):756, 2002. 18

[45] Filip Radenovic, Abhimanyu Dubey, Abhishek Kadian, Todor Mihaylov, Simon Vandenhende, Yash Patel, Yi Wen, Vignesh Ramanathan, and Dhruv Mahajan. Filtering, Distillation, and Hard Negatives for Vision-Language Pre-Training, 2023. arXiv:2301.02280 [cs]. 1, 12

[46] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning Transferable Visual Models From Natural Language Supervision, 2021. Number: arXiv:2103.00020 arXiv:2103.00020 [cs]. 1, 2, 7, 12, 15

[47] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning Transferable Visual Models From Natural Language Supervision, 2021. arXiv:2103.00020 [cs]. 13

[48] Yongming Rao, Wenliang Zhao, Guangyi Chen, Yansong Tang, Zheng Zhu, Guan Huang, Jie Zhou, and Jiwen Lu. Denseclip: Language-guided dense prediction with context-aware prompting. *CoRR*, abs/2112.01518, 2021. 6

[49] Eleanor Rosch. Principles of Categorization. In *Readings in Cognitive Science, a Perspective From Psychology and Artificial Intelligence*, pages 312–22. Morgan Kaufmann Publishers, 1988. 2

[50] Tong Shao, Zhuotao Tian, Hang Zhao, and Jingyong Su. Explore the Potential of CLIP for Training-Free Open Vocabulary Semantic Segmentation, 2024. arXiv:2407.08268 [cs]. 2, 12

[51] Sheng Shen, Liunian Harold Li, Hao Tan, Mohit Bansal, Anna Rohrbach, Kai-Wei Chang, Zhewei Yao, and Kurt Keutzer. HOW MUCH CAN CLIP BENEFIT VISION-AND- LANGUAGE TASKS? 2022. 1

[52] Jaisidh Singh, Ishaan Shrivastava, Mayank Vatsa, Richa Singh, and Aparna Bharati. Learn "No" to Say "Yes" Better: Improving Vision-Language Models via Negations, 2024. arXiv:2403.20312 [cs]. 1, 2, 12

[53] Harald Steck, Chaitanya Ekanadham, and Nathan Kallus. Is Cosine-Similarity of Embeddings Really About Similarity? In *Companion Proceedings of the ACM Web Conference 2024*, pages 887–890, 2024. arXiv:2403.05440 [cs]. 2

[54] Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for reasoning, 2023. 2, 12

[55] Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide shut? exploring the visual shortcomings of multimodal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9568–9578, 2024. 2, 12

[56] Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes Wide Shut? Exploring the Visual Shortcomings of Multimodal LLMs. In *2024 IEEE/CVF*

*Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9568–9578, Seattle, WA, USA, 2024. IEEE. 1, 2

[57] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 6

[58] Feng Wang, Jieru Mei, and Alan Yuille. SCLIP: Rethinking Self-Attention for Dense Vision-Language Inference, 2023. arXiv:2312.01597 version: 1. 2, 12

[59] Haonan Wang, Minbin Huang, Runhui Huang, Lanqing Hong, Hang Xu, Tianyang Hu, Xiaodan Liang, Zhenguo Li, Hong Cheng, and Kenji Kawaguchi. Boosting Visual-Language Models by Exploiting Hard Samples, 2024. arXiv:2305.05208 [cs]. 1, 12

[60] Haicheng Wang, Chen Ju, Weixiong Lin, Shuai Xiao, Mengting Chen, Yixuan Huang, Chang Liu, Mingshuai Yao, Jinsong Lan, Ying Chen, Qingwen Liu, and Yanfeng Wang. Advancing Myopia To Holism: Fully Contrastive Language-Image Pre-training, 2024. arXiv:2412.00440 [cs]. 1, 12

[61] Jeremy M. Wolfe and Todd S. Horowitz. Five factors that guide attention in visual search. *Nature Human Behaviour*, 1(3), 2017. 2

[62] A. D. Wyner. Random packings and coverings of the unit n-sphere. *The Bell System Technical Journal*, 46(9):2111–2118, 1967. 18

[63] Yutaro Yamada, Yingtian Tang, Yoyo Zhang, and Ilker Yildirim. When are Lemons Purple? The Concept Association Bias of Vision-Language Models, 2024. arXiv:2212.12043 [cs] version: 2. 7

[64] Lewei Yao, Runhui Huang, Lu Hou, Guansong Lu, Minzhe Niu, Hang Xu, Xiaodan Liang, Zhenguo Li, Xin Jiang, and Chunjing Xu. Filip: Fine-grained interactive language-image pre-training. In *International Conference on Learning Representations*, 2022. 2, 12

[65] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *Transactions on Machine Learning Research*, 2022. 2, 7, 12

[66] Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. WHEN AND WHY VISION-LANGUAGE MODELS BE- HAVE LIKE BAGS-OF-WORDS, AND WHAT TO DO ABOUT IT? 2023. 1, 2, 7, 8, 12

[67] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. *CoRR*, abs/1905.04899, 2019. 13

[68] Yuhang Zang, Hanlin Goh, Josh Susskind, and Chen Huang. Overcoming the Pitfalls of Vision-Language Model Finetuning for OOD Generalization, 2024. arXiv:2401.15914 [cs]. 12

[69] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid Loss for Language Image Pre-Training, 2023. arXiv:2303.15343 [cs]. 2

[70] Yuhui Zhang, Alyssa Unell, Xiaohan Wang, Dhruba Ghosh, Yuchang Su, Ludwig Schmidt, and Serena Yeung-Levy. Why are Visually-Grounded Language Models Bad at Image Classification?, 2024. arXiv:2405.18415 [cs]. 14

[71] Zheyuan Zhang, Fengyuan Hu, Jayjun Lee, Freda Shi, Parisa Kordjamshidi, Joyce Chai, and Ziqiao Ma. Do Vision-Language Models Represent Space and How? Evaluating Spatial Frame of Reference Under Ambiguities, 2024. arXiv:2410.17385 [cs]. 12

[72] Tiancheng Zhao, Tianqi Zhang, Mingwei Zhu, Haozhan Shen, Kyusong Lee, Xiaopeng Lu, and Jianwei Yin. VL-CheckList: Evaluating Pre-trained Vision-Language Models with Objects, Attributes and Relations, 2023. arXiv:2207.00221 [cs]. 2, 8, 12, 19

[73] Chong Zhou, Chen Change Loy, and Bo Dai. Extract Free Dense Labels from CLIP, 2022. arXiv:2112.01071 [cs]. 6

# Is CLIP ideal? No. Can we fix it? Yes!
# –Supplementary Material–

## A. Extended Related Work

**Vision Language Models.** Recent advancements in VLMs have significantly bridged the gap between vision and language. A seminal work in this area is CLIP [46], which demonstrated that large-scale contrastive learning can effectively align image and text representations into a shared embedding space. Building upon its core principles of large-scale contrastive pretraining and joint representation learning, several subsequent works have explored alternative architectures and training paradigms, including ALIGN [16], FILIP [64], SLIP [36], ALBEF [24], and CoCa [65], to name a few. Autoregressive VLMs [31, 32] have emerged as compelling alternatives to CLIP by jointly attending to both text and image embeddings. Neurosymbolic program synthesis methods like ViperGPT [54], VisProg [14], and VADAR [34] also mitigate some of CLIP's limitations by formulating the text-image semantic distance acquisition as several subproblems. While these methods are more comprehensive than CLIP and can specialize in complex visual reasoning, they are orders of magnitude more expensive to infer and do not offer the same simplicity or gradient retention as CLIP, limiting their downstream applicability. In fact, most of these models rely on a CLIP-like model's latent space as a submodule. As such there remains a strong motivation to continue refining and extending CLIP-like architectures by addressing their inherent shortcomings.

**Empirical Limitations of CLIP.** A growing body of work has revealed several limitations of CLIP in handling complex visual-text interactions. One major issue is its difficulty in distinguishing between different attribute-concept bindings in multi-object scenes [7, 22, 37]. For example, the text prompt "A purple sphere" will have very high cosine similarity with an image that contains a purple cube and a yellow sphere, despite the yellow attribute not belonging to the spherical object. Lewis *et al.* propose CLEVR-bind as a simple benchmark which isolates the attribute binding capability of VLMs. More comprehensive natural-language benchmarks targeting attribute binding include Attribute, Relation, and Order (ARO), Sugarcrepe [15], VL-checklist [72], and Multimodal Visual Patterns (MMVP) [55]. Additional studies have noted that CLIP's text embeddings often behave like "bag-of-words" in practice, leading to imsinterpretations of object layouts or conflates multiple entities within a single scene [22, 37, 66]. Yuksekgonul *et. al* specifically propose WhatsUP as a benchmark which isolates spatial reasoning capacities of VLMs. In addition, the suite of compositional understanding benchmarks (ARO, VL-checklist, Sugarcrepe, MMVP) also include captions that require spatial reasoning. Another notable failure mode is CLIP's inability to accurately represent negation [1, 52]. In response, Alhamoud *et al.* develop NegBench to specifically assess how well VLMs handle various forms of negatory sentences. The aforementioned compositional benchmarks further challenge models with captions that require proper negation understanding.

Some other criticisms of CLIP are its inability to generalize to different reference frames [71] or to count [40]. While these issues represent additional challenges for CLIP, they are beyond the scope of our current work.

**Proposed Solutions to CLIP Limitations.** The most prominent method of corrections have been to change the training data distribution, such as retraining or fine-tuning CLIP with hard negative or positive examples [1, 20, 40, 45, 52, 59, 66] or increasing the training data for more comprehensive and longer captions [28, 60]. However, simply scaling the training data to mitigate specific problems will often lead to reduced generalization [15, 33, 68].

Some engineering solutions include using object detectors to segment images into smaller ROIs [26], adding attribution tracing for correct text-image pairs during training [27], patch clustering for semantic segmentation [50], using chain of thought spatial reasoning [8] or altering the self attention mechanism in the vision encoder [58]. Simpler solutions may be querying CLIP with multiple descriptions [35]. All of these methods require retraining CLIP from scratch, or typically adding a heavy-handed component to single out the objects in a scene.

## B. Experimental Details

### B.1. Implementation Details

By convention the rows of the dense map correspond to one text token embedding, and the columns to one image patch. Every text and image pair creates a DCSM of shape (30,197), where 30 is the maximum number of text tokens and 197 is the number of 16x16 image patches in an image of shape 224x224. Text prompts shorter than 30 tokens are padded with EOS tokens. FRs are therefore of the shape (num-image-patches, 1). For example, for the sentence "An image of a circle above a triangle", the word *above* is a functional word and the corresponding row in the DCSM gets replaced with the respective FR in the lookup table. For synonymous functional words, we use a single FR. (Somewhat unusually, we consider "front","below", and "behind","above" to be synonymous, as DCSMs use a 2D fixed frame of reference due to the topology being represented by the patch index.) The DCSMs are z-score normalized for stable training.

For all training and experiments, we use a lightweight CNN with 2 convolutional layers and a hidden dimension of 128. In addition to a 20-fold reduction in parameters, we train our network with a batch size of 8, a 4000-fold decrease from the original mini-batch size of 32,768. In fact, the sum of all our training data per model is smaller than this number.

Our model outputs a single score for each image and text pair DCSM. During training, we use a contrastive cross-entropy loss as with the original CLIP [47]. We train with the Adam optimizer with learning rate initialized at $1e - 3$. One model is trained with a curated synthetic dataset and another with COCO 2017 training split. Dataset curation is detailed below.

### B.2. Datasets

We train our pipeline on two different datasets - one synthetic dataset composed of open-source 3D assets from Objaverse [12] placed upon randomized backgrounds, and another generated from COCO-train-2017 [30].

We are mainly interested in labeling images with text prompts that lie in the Condition 2,3,4 category. That is, for both the synthetic case and the COCO-train case, we generate a dataset for attribute binding, spatial relationships/localization, and negation. Samples from each dataset are shown in Fig. 6.

**Attribute Binding Dataset** Every image in this dataset includes two distinct objects of unique colors or sizes. For each image, we generate a "hard negative". So if there is an image with a "red cow and purple ghost", we also generate an image with "red ghost and purple cow". This means that every sample has a positive and negative image, and two positive and negative captions each. The positive image

contains object $A$ with attribute $A_{att}$, and object $B$ with attribute $B_{att}$. The negative image contains the same objects but with swapped attributes. The positive caption options are: (P1) "$A_{att}$ $A$ and $B_{att}$ $B$" and (P2) "$B_{att}$ $B$ and $A_{att}$ $A$". The negative caption options are: (N1) "$A_{att}$ $B$ and $B_{att}$ $A$" and (N2) "$B_{att}$ $A$ and $A_{att}$ $B$". We generate 5,402 images for this dataset, which makes 2701 samples with associated opposites.

For the COCO-train set, we use a natural language processing library to extract adjective-noun pairs in the natural language captions, and select images that have at least two distinct objects $A, B$ with distinct attributes $A_{att}, B_{att}$. Captions follow the same format as above. We select 8,547 images from COCO-train towards this dataset.

**Spatial Relationships and Localization** Similarly as above, we generate synthetic images where one object is placed either above, below, to the left, or to the right of, another object with random jitter. For every positive image where $A$ is $rel$ to $B$, there is a negative image where $A$ is $rel_{opp}$ to $A$. Here, (above, below) are opposite relation pairs, as are (left of, right of). The positive caption options are: (P1) "$A$ $rel$ $B$" and (P2) "$B$ $rel_{opp}$ $A$". The negative caption options are: (N1) "$A$ $rel_{opp}$ $B$" and (N2) "$B$ $rel$ $A$". We generate 11,324 images for this dataset, which makes 5662 samples with associated opposites.

For the COCO images, we choose images where at least two distinct objects are present, and use the relationship between their bounding boxes to validate that they satisfy the definition of one of the spatial relationships being considered. Note that, as we are using patch location on the image to preserve topology, we use a fixed frame of reference to determine the meaning of "above", "below", "left", and "right". To make a negative version of the image, we use the CutMix technique [67] to swap the image content in the two bounding boxes. With the positive images from COCO-train and generated hard negatives, there are 11,502 images in this dataset, which makes 5751 samples with associated opposites.

**Negaton** Generating negatory captions is tricky. With contrastive training, an entirely negatory caption (e.g., "An image without a turtle") necessitates that all other images in the batch be a positive sample (e.g., they must now all contain a turtle).

As such, for the first pass for the models trained with synthetic data, we generate images with two random objects $A1$ and $A2$, and choose as its hard negative another image which contains two non-overlapping objects, $B1$ and $B2$. For the negatory term $\neg$, we choose between ['but not', 'and no', 'without']. The positive caption options are: (P1) "$A1\neg B1$" and (P2) "$A2\neg B2$". The negative caption options are: (N1) "$B1\neg A1$" and (N2) "$B2\neg A2$". We generate 10,000 images for this dataset.

For COCO-train we select two images that have two dis-

Figure 6. Overview of Dataset Curation. Each box is a dataset. The left image is a positive pair with the top two captions, and the right image is a positive pair with the bottom two captions. The COCO-train Attribute set does not have a hard negative image counterpart.

tinct object labels, and generate captions the same way. We select 10,000 images toward this dataset.

### B.3. Pseudocode

Below we illustrate the process for extracting DCSMs.

```python
dense_image_features = image_encoder(I).
    last_hidden_state.unsqueeze(1)
# Shape: (batch_size, 1, iseq, embed_dim)

dense_text_features = text_encoder(T).
    last_hidden_state.unsqueeze(0)
# Shape: (1, batch_size, tseq, embed_dim)

dcsm = einsum( "bqie, lpte -> bpit",
    dense_image_features, dense_text_features )
# Shape: (batch_size, batch_size, iseq, tseq)

dcsm = add_functional_rows(dcsm, lookup_list)

out = lightweight_cnn(dcsm)
# Shape: batch_size, batch_size

labels = eye(out.shape[0])
# Shape: batch_size, batch_size

loss = CE(out, labels) + CE(out.t(), labels.t())
```

Figure 7. Pytorch-like code for training our model with DCSMs

## C. Ablation Studies

### C.1. Impact of DCSMs

In Sec 5. we said that it follows from our analysis that neither a fine-tuned/reprojected CLIP embedding space, nor a learned scoring module, could alone be the fix to CLIP's

| Model | WhatsUP | COCO$_{2obj}$ |
|---|---|---|
| **Ours - DCSM (CNN)** | **62.6** | **70.9** |
| **Ours - DCSM (CNN)**$_{w/oFR}$ | 48.6 | 55.5 |
| **Ours - DCSM (ViT)** | 29.4 | 47.0 |
| CLIP - ViTB/16 | 30.5 | 45.9 |
| CLIP - ViTB/16$_{f.t.\ synth}$ | 25.5 | 49.1 |
| CLIP - ViTB/16$_{MLP\ scorer}$ | 25.4 | 53.6 |

Table 4. The DCSM networks were trained with synthetic data.

fundamental shortcomings. To verify this conclusion, we perform a series of ablations. First, we finetune OpenAI's CLIP (ViTB/16) with the same synthetic dataset used for our model training. We also train a MLP scoring module which takes concatenated text and image embeddings as input to output a score, again on the same synthetic dataset. Both attempts fail to improve performance on WhatsUP and COCO-spatial, resulting in accuracies near chance (25% and 50%, respectively).

Further, we alter our training pipeline in two different ways to assess the need for a CNN as well as the FRs. Removal of the FRs decreases performance overall, but the increased information capacity from using the DCSMs and the downstream network still allows the model to perform above fine-tuned CLIP models. Replacement of the CNN with a comparably small ViT, with a patch size of 2 and 2 layers of 4 attention heads, resulted in another near-chance performance on the datasets. The ViT appears more prone to overfitting to the training set, as it does not have the imposed constraint of pattern-recognizing kernels as in CNNs.

All networks were trained with learning rate 1e-3 with Adam for 27 epochs. Under minimal compute, the CNN generalized much better. Fine-tuning CLIP projection layers with the same dataset for the same number of epochs did not result in any noticeable performance increase.

### C.2. Classification Performance

A known problem with VLMs using CLIP embeddings is the decline in classification capacity [70]. We evaluate our model on two classification datasets and find that, despite being trained with simple prompts and no image classification captions, the reduction in classification performance is minimal compared to those observed in BLIP or LLaVA.

### C.3. Scaling Analysis

In this work, we showcase a very small and computationally light network and training pipeline for the DCSM method. To verify that this method will scale with increasing data, we perform a scaling analysis. Fig. 8 shows the results. The x axis is the approximate number of samples from the curated COCO2017 training set. From this we see that our training pipeline is likely to scale with increasing the dataset.

| Model | Caltech101 | Flowers102 |
|---|---|---|
| CLIP-ViT B/16 [46] | **82.6** | **67.7** |
| **Ours - DCSM + synth** | **77.8** | **52.9** |
| **Ours - DCSM + coco** | **79.2** | 43.7 |
| BLIP2-2.7B [25] | 22.3 | 14.2 |
| IBLIP-7B [10] | 58.4 | 26.8 |
| LLaVA1.5-7B [32] | 62.1 | 10.2 |

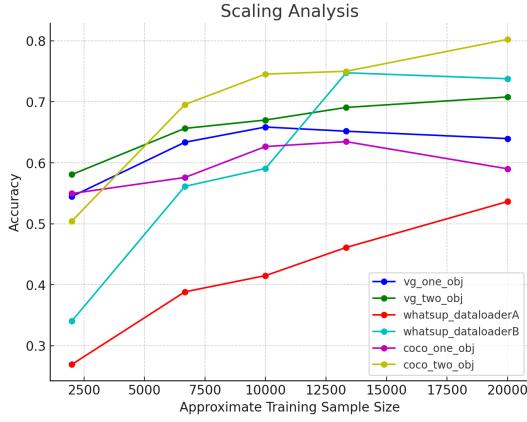Table 5. Classification Accuracy of VLMs. Top three scores per dataset are bolded.



Figure 8. Result of linearly scaling training data. x values are approximate dataset sizes.

# D. Empirical Observations of CLIP shortcomings

In Fig. 9 we show the empirical effects of the superposition derived in Lemma 1. In summary, the figure serves to show how an image with an increasingly greater number of objects present embeds increasingly farther from the text label for any one of those objects in the CLIP latent space. The degree of this effect is such that beyond 6-8 objects in one image, CLIP embeddings of random noise images are similarly close to those object text labels.

# E. Definitions and Proofs

## E.1. Definitions Addendum

**Definition 1.** For x = the concept of a bird, its mapping in $\mathbb{I}$ is an image of a bird, and its mapping in $\mathbb{T}$ is "bird". For all concepts x in $\mathbb{V}$, there is an equivalence class in $\mathbb{I}$ and $\mathbb{T}$ that unambiguously represent that concept, respectively. We narrow the scope of all images to 1 representative image per concept. So for the object concept of a "bird", there is one image which is equivalent. Similarly, for each concept x there is only 1 corresponding $d \in \mathbb{T}$.

We also choose $M$ object concept elements toward the subset $\mathbb{V}$. In the manner of large ontologies (such as Im-
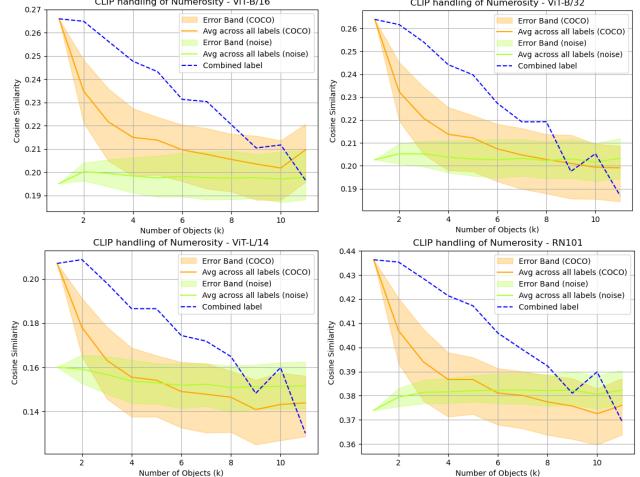


Figure 9. For each image in COCO-validation, we identify k objects with labels. Then we take the cosine similarity between all k labels and the image. The orange line shows the average cosine similarity for images with k objects and all labels that appear in the image. The blue dotted line shows the cosine similarity for a label which combines all k labels and the image being considered, averaged for each k. The green plot shows the average cosine similarity score between 5 random noise images and the labels of all COCO-validation images. The error bands indicate the 25th and 75th quartile.

ageNet1000), we take it to be true that there exist at least $M \leq 1000$ object classes which are mutually exclusive. This means that for any representative image in the set of $M$, an expert human reviewer would be able to assign it to that unique object class. These object classes must be *subordinate categories* on any hierarchical directed tree of visual concepts. (An example concept hierarchy: *car* and *bus* are lower-level concepts that are below *vehicle*.) This ensures no two concepts in $\mathbb{V}$ overlap in semantics.

## E.2. Conditions Addendum

**Condition 1.** Semantic separability means: if human annotators would agree that a text caption appropriately describes an image, the embeddings for that caption and image should have high cosine similarity. Conversely, the score should be low if an annotator deems the caption to be inappropriately matched to the image.

Notably, this condition does not require correct hierarchical organization among concepts or semantically accurate placement of synonymous phrases. In reality there exist more than N distinct concept categories and each category may have synonymous text representations and multiple instances and viewpoints of the representative visual object. Our goal is to set up extremely minimal geometric requirements for $C$ to evaluate whether they are possible to attain.

This condition must be satisfied for $C$ to perform zero-

shot image classification, retrieval, and semantic similarity search. Example benchmarks that pose this challenge include Imagenet, COCO, LAION, and many more.

**Condition 2.** This condition must be satisfied for $C$ to perform tasks where the model must identify attributes associated with different objects in a scene. This could be towards scene understanding, vision question answering, or accurate image retrieval. Specific datasets include CLEVR-bind, NCDataset-grayscale, VL-checklist, Sugarcrepe, ARO, and MMVP.

**Condition 3.** This condition must be satisfied for $C$ to to perform tasks that require compositional image understanding. This could be for image captioning, text-guided image generation, spatial navigation, and more. Specific datasets include WhatsUP, Coco-spatial, MMVP, etc.

**Condition 4.** This condition must be satisfied for $C$ to perform well on vision-language tasks that include prompts with negations. Note that this condition is a very relaxed interpretation of negation: Strictly semantically speaking, "not X" is a correct semantic pair with any image that does not have X, requiring a cosine similarity near 1. But the definition of the negation condition we impose does not require such granularity. Again, we seek to pose minimal constraints to identify whether there is some version of negation CLIP could attain under ideal settings. Specific datasets that require negation understanding include NegBench, as well as other compositional benchmarks like VLM-checklist, Sugarcrepe, or MMVP.

### E.3. Lemma 2. Addendum

**Derivation of** $p$**.** We define the perturbed vector:
$$\mathbf{i}(x_a) = (1 - \delta)\,\mathbf{i}(x) + p\mathbf{t}(a), \|\mathbf{i}(x_a)\|^2 = 1 \quad (10)$$
Expanding the norm, we have:
$$\|\mathbf{i}(x_a)\|^2 = \|(1 - \delta)\mathbf{i}(x) + p\mathbf{t}(a)\|^2$$
$$\|\mathbf{i}(x_a)\|^2 = (1 - \delta)^2\|\mathbf{i}(x)\|^2 + p^2\|\mathbf{t}(a)\|^2 \quad (11)$$
$$+ 2(1 - \delta)p\mathbf{i}(x) \cdot \mathbf{t}(a)$$
Since $\mathbf{i}(x)$ and $\mathbf{t}(a)$ are unit vectors, this simplifies to:
$$\|\mathbf{i}(x_a)\|^2 = (1 - \delta)^2 + p^2 + 2(1 - \delta)p\cos\theta \quad (12)$$
where $\cos\theta = \mathbf{i}(x) \cdot \mathbf{t}(a)$. For $\mathbf{i}(x_a)$ to be a unit vector, we set the right hand side to 1:
$$(1 - \delta)^2 + p^2 + 2(1 - \delta)p\cos\theta = 1$$
$$-2\delta + \delta^2 + p^2 + 2(1 - \delta)p\cos\theta = 0 \quad (13)$$
Notice that this is now a quadratic equation in $p$:
$$p^2 + 2(1 - \delta)p\cos\theta - (2\delta - \delta^2) = 0. \quad (14)$$
Use the quadratic formula:
$$p = \frac{-2(1-\delta)\cos\theta \pm \sqrt{[2(1-\delta)\cos\theta]^2 + 4(2\delta - \delta^2)}}{2} \quad (15)$$
We can simplify the above to find the correct value of $p$ that ensures $\mathbf{i}(x_a)$ is a unit vector.
$$p = -(1-\delta)\cos\theta \pm \frac{\sqrt{4(1-\delta)^2\cos^2\theta + 8\delta - 4\delta^2}}{2}. \quad (16)$$

**Analysis for noise vectors.** In Lemma 2 we derived

$$\boxed{\begin{aligned}\mathbf{i}(x_a, y_b) &= \frac{(1 - \delta)(\mathbf{i}(x) + \mathbf{i}(y)) + p\mathbf{t}(a) + q\mathbf{t}(b)}{2} \\ &= \mathbf{i}(x_b, y_a)\end{aligned}} \quad (17)$$

using $\mathbf{i}(x_a) = (1 - \delta)\mathbf{i}(x) + p\mathbf{t}(a)$.

Now we show that analytically, the inclusion of a noise vector $\epsilon$ does not change the results. Specifically, we want to simulate
$$\mathbf{i}(x_a) = \frac{(1 - \delta)\mathbf{i}(x) + \delta\mathbf{t}(a) + \epsilon}{\|norm\|}$$
for some randomly sampled $\epsilon$ and $\delta$.

In Fig. 10 we showcase the results of sampling $\epsilon$ from a standard normal distribution, with varying weights.
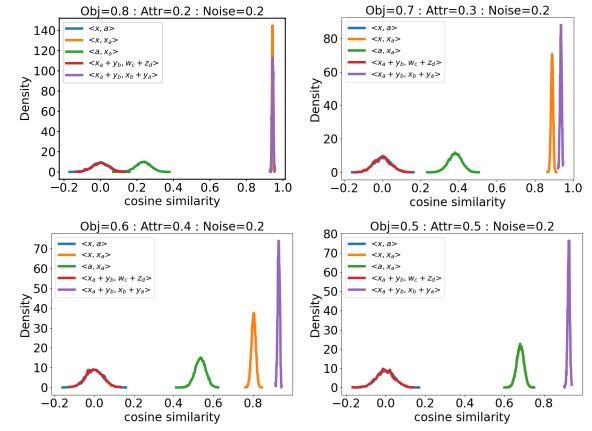


Figure 10. Titles of each subplot indicate the weights of the different components composing $\mathbf{i}(x_a, y_b), \mathbf{i}(x_b, y_a)$. Obj= indicates weight of the object concept embeddings $\mathbf{i}(x), \mathbf{i}(y)$, Attr= indicates weight of the attribute text embeddings $\mathbf{t}(a), \mathbf{t}(b)$, and Noise= indicates weight of the noise vector $\epsilon$. In the legend, object concepts and attribute embeddings are denoted in shorthand: $x = \mathbf{i}(x), a = \mathbf{t}(a)$, and so on.

We observe the following relations remain consistent:
- As expected of randomly initialized vectors in high dimensions: $\mathbf{i}(x) \cdot \mathbf{t}(a) \approx 0$
- For some unrelated object-attribute pairs, their image embeddings are roughly orthogonal as well: $\mathbf{i}(w_c, z_d) \cdot \mathbf{i}(x_a, y_b) \approx 0$
- The strong conclusion from Lemma 2 is approximately always true: $\mathbf{i}(x_b, y_a) \cdot \mathbf{i}(x_a, y_b) \approx 1$

regardless of $\mathbf{i}(x_a) \cdot \mathbf{t}(a), \mathbf{i}(x) \cdot \mathbf{i}(x_a)$. As such, we see that even if there is noise in the superpositions described in Lemma 1 and 2, $C$ still cannot disambiguate between different pairings of the same two attributes and two objects.

### E.4. Contradiction for Condition 1 and 3

Now we show Condition 3 cannot be met if Condition 1 is met. Below we will show two impossibility cases and prove

them.

**Lemma 3.** $C$ **cannot accurately represent both the distance between spatial locations and relationships at the same time.**

*Proof.* We first derive $\mathbf{i}(x, g^{<rel>}, y)$ for two antonymous $g^{<rel>}$ to satisfy Condition 1.2. Then we consider an example scenario with two objects, two spatial relationships, and two localizing terms. We will find that for three sample images, the cosine similarities between their embeddings and four textual clauses will have to contradict Condition 3 for some pairs.

For two concepts $x$ and $y$, their combined embedding is Eq. (3). Now, if we want to express some compositional relationship $g^{<rel>} \in \mathbb{G}$ between x and y such that $\mathbf{i}(x, g_1^{<rel>}, y) \neq \mathbf{i}(x, g_2^{<rel>}, y)$, we can write

$$\mathbf{i}(x, g_1^{<rel>}, y) = (1 - \delta)\,\mathbf{i}(x, y) \;+\; \mathbf{v}_1$$
$$\mathbf{i}(x, g_2^{<rel>}, y) = (1 - \delta)\,\mathbf{i}(x, y) \;+\; \mathbf{v}_2 \qquad (18)$$

where $\delta \ll 1$ and $\mathbf{v}_1 \neq \mathbf{v}_2$. Similar to Lemma 2, $\mathbf{v}$ is a small location-specific component, as $\mathbf{i}(x, g^{<rel>}, y)$ must remain close to $\mathbf{i}(x, y)$ per Condition 1.2.

Let $g_L^{<rel>} = g_L$ be the relational concept whose equivalent mapping in $\mathbb{T}$ is "_ left of _", and $g_R^{<rel>} = g_R$ "_ right of _". Then we can write:

$$\mathbf{i}(x, g_L, y) = (1 - \delta)\,\mathbf{i}(x, y) + \mathbf{e}_{\perp,L}$$
$$\mathbf{i}(x, g_R, y) = (1 - \delta)\,\mathbf{i}(x, y) + \mathbf{e}_{\perp,R}$$

where $\mathbf{e}_{\perp,L}, \mathbf{e}_{\perp,R}$ both lie in the orthogonal error subspace (of dimension $N - 1$) and have fixed magnitude $\sqrt{2\delta - \delta^2}$ such that $\mathbf{i}(x, g_L, y)$ is a unit vector. Notice that we cannot use the intuition from Lemma 2 that $v$ must be composed of the textual component of $g_L$ - an image with a melon above a bed does not intuitively need to embed closely with the text embedding for "above".

Then to *maximize* the Euclidean distance $\|\mathbf{i}(x, g_L, y) - \mathbf{i}(x, g_R, y)\|$, we must choose $\mathbf{e}_{\perp,R} = -\mathbf{e}_{\perp,L}$. Since $(x, g_L, y), (x, g_R, y) \in \mathbb{S}$, they also have an equivalent item in $\mathbb{T}$.
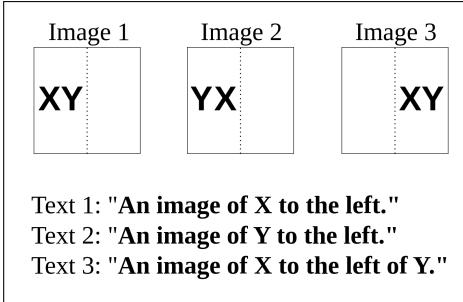


| Image 1 | Image 2 | Image 3 |
| :---: | :---: | :---: |
| **XY** | **YX** | **XY** |

Text 1: **"An image of X to the left."**
Text 2: **"An image of Y to the left."**
Text 3: **"An image of X to the left of Y."**

Figure 11. Simple setup. Here the only relations that exist are "_ left _", "_ right of _", "_ to the left", and "_ to the right".

Now we formulate a proof by contradiction. Consider three images and three text prompts, as shown in

Fig. 11. In addition to $g_L, g_R$, we also introduce $g_l^{<loc>} = g_l, g_r^{<loc>} = g_r$, to represent "_ to the left" and "_ to the right" in $\mathbb{T}$, respectively. These are the only four compositional concepts we consider for simplicity, but conclusions will generalize.

We denote $(1 - \delta)\,\mathbf{i}(x, y) = \mathbf{e}_{||}$, the three images as $\mathbf{i}(\text{image}_k) = \mathbf{i}_k$ and $\mathbf{t}(\text{text}_j) = \mathbf{t}_j$, for $k, j \in 1, 2, 3$. $\mathbf{i}_k = \mathbf{e}_{||} + \mathbf{e}_{\perp,k}$ where $\mathbf{e}_{||}$ is the shared parallel component for all $\mathbf{i}_k$.

If it is possible to sufficiently represent Condition 3.1 and 3.2 in any $C$, then there must exist a $C$ in which $\mathbf{e}_{\perp,1} \cdot \mathbf{e}_{\perp,2} > 0$ and $\mathbf{e}_{\perp,1} \cdot \mathbf{e}_{\perp,3} > 0$. This is because $\mathbf{e}_{\perp,1} \cdot \mathbf{e}_{\perp,2} = 1$ indicates the two error vectors represent synonymous compositional semantics, $\mathbf{e}_{\perp,1} \cdot \mathbf{e}_{\perp,2} = -1$ indicates antonymous semantics, and error vectors with irrelevant semantics are orthogonal to each other. Therefore, as images where X and Y are both to the left of the image must have higher cosine similarity with each other than with other images containing x and y, $\mathbf{e}_{\perp,1} \cdot \mathbf{e}_{\perp,2} > 0$ must be true. Similarly, since images where X is to the left of Y should have higher cosine similarity with each other, $\mathbf{e}_{\perp,1} \cdot \mathbf{e}_{\perp,3} > 0$ must be true.

For each image in Fig. 11, the embeddings must optimize the following local similarities:

$$\mathbf{i}_1 = \underset{\mathbf{i}_1}{\arg\max}(\mathbf{i}_1 \cdot \mathbf{t}_1 + \mathbf{i}_1 \cdot \mathbf{t}_2 + \mathbf{i}_1 \cdot \mathbf{t}_3)$$
$$\mathbf{i}_2 = \underset{\mathbf{i}_2}{\arg\max}(\mathbf{i}_2 \cdot \mathbf{t}_1 + \mathbf{i}_2 \cdot \mathbf{t}_2 + \mathbf{i}_2 \cdot -\mathbf{t}_3) \qquad (19)$$
$$\mathbf{i}_3 = \underset{\mathbf{i}_3}{\arg\max}(\mathbf{i}_3 \cdot -\mathbf{t}_1 + \mathbf{i}_3 \cdot -\mathbf{t}_2 + \mathbf{i}_3 \cdot \mathbf{t}_3)$$

subject to $\|\mathbf{i}_k\| = 1$. This allows us to solve for $\mathbf{e}_{\perp,k}$s. For $k = 1$:

$$\mathbf{e}_{\perp,1} = \underset{\mathbf{e}_{\perp,1}}{\arg\max}((\mathbf{e}_{||} + \mathbf{e}_{\perp,1}) \cdot \mathbf{t}_1 + ((\mathbf{e}_{||} + \mathbf{e}_{\perp,1}) \cdot \mathbf{t}_2 +$$
$$((\mathbf{e}_{||} + \mathbf{e}_{\perp,1}) \cdot \mathbf{t}_3) \qquad (20)$$

subject to $\|\mathbf{e}_{\perp,1}\| = \sqrt{2\delta - \delta^2}$. $\mathbf{e}_{||}$ is fixed so this becomes:

$$\mathbf{e}_{\perp,1} = \underset{\mathbf{e}_{\perp,1}}{\arg\max}(\mathbf{e}_{||} \cdot (\mathbf{t}_1 + \mathbf{t}_2 + \mathbf{t}_3) + \mathbf{e}_{\perp,1} \cdot (\mathbf{t}_1 + \mathbf{t}_2 + \mathbf{t}_3))$$
$$= (\mathbf{t}_1 + \mathbf{t}_2 + \mathbf{t}_3) \cdot \sqrt{2\delta - \delta^2} \qquad (21)$$

Similarly, we get:

$$\mathbf{e}_{\perp,2} = (\mathbf{t}_1 + \mathbf{t}_2 - \mathbf{t}_3) \cdot \sqrt{2\delta - \delta^2}$$
$$\mathbf{e}_{\perp,3} = (-\mathbf{t}_1 - \mathbf{t}_2 + \mathbf{t}_3) \cdot \sqrt{2\delta - \delta^2} \qquad (22)$$

Now taking the dot products, we have:

$$\mathbf{e}_{\perp,1} \cdot \mathbf{e}_{\perp,2} = (|\mathbf{t}_1| + |\mathbf{t}_2| - |\mathbf{t}_3|) \cdot 2\delta = 2\delta$$
$$\mathbf{e}_{\perp,1} \cdot \mathbf{e}_{\perp,3} = (-|\mathbf{t}_1| - |\mathbf{t}_2| + |\mathbf{t}_3|) \cdot 2\delta = -2\delta \qquad (23)$$
$$\mathbf{e}_{\perp,2} \cdot \mathbf{e}_{\perp,3} = (-|\mathbf{t}_1| - |\mathbf{t}_2| - |\mathbf{t}_3|) \cdot 2\delta = -6\delta$$

Practically speaking, it is possible by adding more representative samples in the training dataset to change the weights of $\mathbf{t}_j$s. That is, for some $\beta_1 + \beta_2 + \beta_3 = 3$,

Eqs. (21,22) could be reformulated as:

$$\mathbf{e}_{\perp,1} = (\beta_1\mathbf{t}_1 + \beta_2\mathbf{t}_2 + \beta_3\mathbf{t}_3) \cdot \sqrt{2\delta - \delta^2}$$

$$\mathbf{e}_{\perp,2} = (\beta_1\mathbf{t}_1 + \beta_2\mathbf{t}_2 - \beta_3\mathbf{t}_3) \cdot \sqrt{2\delta - \delta^2} \qquad (24)$$

$$\mathbf{e}_{\perp,3} = (-\beta_1\mathbf{t}_1 - \beta_2\mathbf{t}_2 + \beta_3\mathbf{t}_3) \cdot \sqrt{2\delta - \delta^2}$$

The dot products then become:

$$\boxed{\begin{aligned} \mathbf{e}_{\perp,1} \cdot \mathbf{e}_{\perp,2} &= 2\delta\left(\beta_1^2 + \beta_2^2 - \beta_3^2\right) \\ \mathbf{e}_{\perp,1} \cdot \mathbf{e}_{\perp,3} &= 2\delta\left(-\beta_1^2 - \beta^2 + \beta_3^2\right) \\ \mathbf{e}_{\perp,2} \cdot \mathbf{e}_{\perp,3} &= -2\delta\left(\beta_1^2 + \beta_2^2 + \beta_3^2\right) \end{aligned}} \qquad (25)$$

Note that regardless of the reweighting, $\mathbf{e}_{\perp,1} \cdot \mathbf{e}_{\perp,2} = -\mathbf{e}_{\perp,1} \cdot \mathbf{e}_{\perp,3}$. This directly negates our previous observation that an ideal $C$ must satisfy $\mathbf{e}_{\perp,1} \cdot \mathbf{e}_{\perp,2} > 0$ and $\mathbf{e}_{\perp,1} \cdot \mathbf{e}_{\perp,3} > 0$. As such, there exists no $C$ which sufficiently represents both relational and objective space in the image embeddings.

□

**Lemma 4.** $C$ **cannot accurately represent compositional concepts of different hierarchy.**

*Proof.* Here we will show that general prepositions are erroneously closer to all unrelated prepositions.

Some prepositions are more general than others. For example, $g_B^{<rel>} = g_B$ where $f_{G,T}(g_B) = $ ”_ beside _” semantically includes both $g_L$ and $g_R$. The ideal placement for $\mathbf{t}(x, g_B, y)$ should locally optimize for the following similarities:

$$\mathbf{t}(x, g_B, y) = \underset{\mathbf{t}(x, g_B, y)}{\operatorname{argmax}} \Big[ \mathbf{t}(x, g_L, y) \cdot \mathbf{t}(x, g_B, y)$$

$$+ \mathbf{t}(x, g_R, y) \cdot \mathbf{t}(x, g_B, y) \qquad (26)$$

$$- \sum_{z=1}^{|\mathbb{G}\setminus\{g_L, g_R, g_B\}|} \mathbf{t}(x, g_z, y) \cdot \mathbf{t}(x, g_B, y) \Big]$$

We know from Lemma 1 that this means $\mathbf{t}(x, g_B, y)$ should be a weighted superposition of $\mathbf{t}(x, g_L, y)$ and $\mathbf{t}(x, g_R, y)$. Since we know $\mathbf{e}_L = -\mathbf{e}_R$, we can write:

$$\mathbf{t}(x, g_L, y) = (1 - \delta) \cdot \mathbf{t}(x, y) + \mathbf{e}_L$$

$$\mathbf{t}(x, g_R, y) = (1 - \delta) \cdot \mathbf{t}(x, y) - \mathbf{e}_L$$

Then the superposition of the two becomes $\mathbf{t}(x, y)$. But this is a semantically erroneous placement for $\mathbf{t}(x, g_B, y)$, as it will incorrectly be closer to any other instance of $x$ and $y$ co-appearing in a scene than the average $\mathbf{t}(x, g^{<rel>}, y)$. For example, for $g_A^{<rel>} = g_A$ where $f_{G,T}(g_A) = $ ”_ above _”, $\mathbf{t}(x, g_B, y) \cdot \mathbf{i}(x, g_A, y) > \mathbf{t}(x, g_L, y) \cdot \mathbf{i}(x, g_A, y)$ even though the two captions are equally inapplicable.

□

### E.5. Contradiction for Condition 1 and 4

Now we show Condition 4 cannot be met if Condition 1 is met. Before we move to the proof, we first discuss in greater detail the possible arrangements of object concepts in $C$.

**Orthogonality.** One straightforward intuition is that since $C$ is a high dimensional space, for any two random concepts $x^1, x^2 \in \mathbb{V}$, they should be approximately orthogonal [62]:

$$\mathbf{t}(x^j) \cdot \mathbf{t}(x^k) \approx 0 \quad \forall j \neq k \qquad (27)$$

This makes it trivial to derive that $\mathbf{t}(\neg x) \cdot \mathbf{t}(y) \approx \mathbf{t}(x) \cdot \mathbf{t}(y) = 0$, violating Condition 4.2. In order to be more rigorous, we show that negation cannot be achieved even under strong idealistic conditions, where $\mathbf{t}(x)$ are uniformly distributed and *distinguishable* from one another. This requires perfect isotropy of $M$ concepts.

**Isotropy.** Starting with $|\mathbb{V}| = M \leq N$ concepts, we determine the ideal distribution for $\mathbf{t}(x)$ $\forall x \in \mathbb{V}$. As in Lemma 1, we denote distinct concepts as: $\mathbf{t}(x^1), \mathbf{t}(x^2), ....$. Since all $x \in \mathbb{V}$ are mutually exclusive in semantics by definition, the distance between any two arbitrary concepts $x^1, x^2$ should be comparable to the distance between $x^1, x^k$ for some $x^k \in \mathbb{V}\setminus\{x^2\}$. Then $C$ must minimize the variance among the cosine similarities of all pairs of concepts:

$$\min_{\mathbf{t}(x^j) \cdot \mathbf{t}(x^k)} \sum_{j=1}^{M} \sum_{k>j}^{M} \left(\mathbf{t}(x^j) \cdot \mathbf{t}(x^k) - \bar{s}\right)^2 \qquad (28)$$

$$\text{s. t. } \|\mathbf{t}(x^k)\| = 1, \quad \forall x^k \in \mathbb{V}$$

where $\bar{s}$ is the mean cosine similarity:

$$\bar{s} = \frac{1}{\binom{M}{2}} \sum_{j=1}^{M} \sum_{k>j}^{M} \mathbf{t}(x^j) \cdot \mathbf{t}(x^k) \qquad (29)$$

Let $\mathbf{t}(x^j) \cdot \mathbf{t}(x^k) = s_{jk}$. Then the objective simplifies to:

$$\min_{s_{jk}} \sum_{j=1}^{M} \sum_{k>j}^{M} \left(s_{jk}^2 - \bar{s}^2\right) = \min_{s_{jk}} \sum_{j=1}^{M} \sum_{k>j}^{M} (s_{jk})^2 - \bar{s}^2 \binom{M}{2} \qquad (30)$$

Differentiate the first and second terms with respect to $s_{jk}$:

$$\frac{\partial}{\partial s_{jk}} \sum_{j=1}^{M} \sum_{k>j}^{M} s_{jk}^2 = 2s_{jk}, \qquad (31)$$

$$\frac{\partial}{\partial s_{jk}} \left(\bar{s}^2 \binom{M}{2}\right) = 2\binom{M}{2}\bar{s}\frac{\partial\bar{s}}{\partial s_{jk}} = \frac{1}{\binom{N}{2}}2\binom{M}{2}\bar{s} = 2\bar{s} \qquad (32)$$

This means that the optimum of Eq. (28) is reached when

$$2\bar{s} - 2s_{jk} = 0 \quad \forall j, k \in (1, M), j \neq k \qquad (33)$$

In other words, the cosine similarities between any two object concepts must be the same as the average. That requires for all $\mathbf{t}(x^k)$ to be isotropically distributed in $\mathbb{R}^N$. The optimal arrangement of all $\mathbf{t}(x^k)$ is then a $M$-1 dimensional regular simplex, which is a structure of $M$ equiangular unit vectors in $\mathbb{R}^M$. Then we have that the cosine similarity of two random vectors in $C$ is:

$$\boxed{\mathbf{t}(x^j) \cdot \mathbf{t}(x^k) = -\frac{1}{M-1} \quad \forall j \neq k} \qquad (34)$$

as all vector pairs in a regular simplex have angles $\arccos\left(-\frac{1}{M-1}\right)$ [44].

**Lemma 5.** **Even under isotropic concept distribution,** $C$

**cannot accurately represent negation.**

*Proof.* We first derive what $\mathbf{t}(x), \mathbf{t}(\neg x)$ must be to respect Conditions 4.1 and 4.2. Then we see that this derivation contradicts Condition 4.3.

For $C$ to ideally represent negation, the following must be true:

$$\mathbf{t}(\text{"}\neg x\text{"}) \cdot \mathbf{i}(x) < \mathbf{t}(\text{"}\neg x\text{"}) \cdot \mathbf{i}(v) \quad (1)$$
$$\mathbf{t}(\text{"}\neg x\text{"}) \cdot \mathbf{i}(v) > \mathbf{t}(\text{"}x\text{"}) \cdot \mathbf{i}(v) \quad (2) \qquad (35)$$
$$\mathbf{t}(\text{"}x\text{"}) \cdot \mathbf{t}(\text{"}y\text{"}) < \mathbf{t}(\text{"}\neg x\text{"}) \cdot \mathbf{t}(\text{"}\neg y\text{"}) \quad (3)$$

for all $v \in \mathbb{V}\backslash\{x\}$

To achieve Eq. (35.1), we solve for:

$$\mathbf{t}(\neg x) = \underset{\mathbf{t}(\neg x)}{\operatorname{argmax}} \left[ \sum_{v=1}^{|\mathbb{V}\backslash\{x\}|} \mathbf{t}(\neg x) \cdot \mathbf{i}(v) - \mathbf{t}(\neg x) \cdot \mathbf{i}(x) \right]$$

$$= \underset{\mathbf{t}(\neg x)}{\operatorname{argmax}} \left[ \mathbf{t}(\neg x) \cdot \left( \sum_{v=1}^{|\mathbb{V}|} \mathbf{i}(v) - \mathbf{i}(x) - \mathbf{i}(x) \right) \right]$$

$$(36)$$

Here, $\sum_{v=1}^{|\mathbb{V}|} \mathbf{i}(v)$ is the sum of all vectors in a regular simplex, which is 0. As such, we find that:

$$\mathbf{t}(\neg x)^* = -\mathbf{i}(x) \qquad (37)$$

Recall that for two vectors in an $M-1$ dimensional simplex, $\mathbf{t}(x^j) \cdot \mathbf{t}(x^k) = -\frac{1}{M-1} \ \forall j \neq k$. With the above solution, we now have that:

$$\mathbf{t}(\neg x^j) \cdot \mathbf{t}(x^k) = \frac{1}{M-1} > \mathbf{t}(x^j) \cdot \mathbf{t}(x^k) \qquad (38)$$

which satisfies condition 4.2. But 4.3 fails due to the following:

$$\boxed{\begin{aligned} \mathbf{t}(\neg x^j) \cdot \mathbf{t}(\neg x^k) = \mathbf{t}(x^j) \cdot \mathbf{t}(x^k) = -\frac{1}{M-1} \\ \mathbf{t}(\neg x^j) \cdot \mathbf{t}(x^k) > \mathbf{t}(\neg x^j) \cdot \mathbf{t}(\neg x^k) \end{aligned}} \qquad (39)$$

Let's say $x^j = $ "chair" and $x^k = $ "penguin". The first erroneous semantic relationship that emerges is that the distance between "chair" and "penguin", which are fully contradictory statements, will be equivalent to the distance between "Not chair" and "Not penguin". For the latter two prompts there exist a lot of images that would be a true match for both, whereas for the first prompt there exists only one.

The second erroneous conclusion is that the cosine similarity between "Not chair" and "penguin" is greater than the cosine similarity between "Not chair" and "Not penguin". This is semantically incorrect for the same reason as above. $\square$

## F. Open Vocabulary Experimental Details

We train two model types, one on 5k COCO images from the training split for 12 epochs, and another on 10k COCO images for 6 epochs. The images each have a hard positive

---

**Original Prompts:**

['A baby elephant stands next to its mother.', 'A desktop computer sitting on top of a wooden desk.', 'A mother stands next to its baby elephant.', 'A wooden computer sitting on top of a desk.']

**Original Lookup Table:**

[ 'above', 'below', 'many', 'no', 'small', 'big', 'not', 'without', 'left', 'right', 'absent', 'but', 'large' ]

**Simplified Prompts:**

[ 'Baby elephant next to mother', 'Desktop computer on wooden desk', 'Mother next to baby elephant', 'Wooden computer on desk' ]

**New Lookup Table:**

['above', 'below', 'many', 'no', 'small', 'big', 'not', 'without', 'left', 'right', 'absent', 'but', 'large', 'near', 'on']

**New Prompts:**

['Baby elephant ⟨NEAR⟩ mother', 'Desktop computer ⟨ON⟩ wooden desk', 'Mother ⟨NEAR⟩ baby elephant', 'Wooden computer ⟨ON⟩ desk']

Table 6. Examples of LLM-in-the-loop natural language simplification and FR extraction.

and negative, where the latter is the same as the former save for two nouns being swapped. We choose this particular type of intervention as CLIP-like VLM performance across the board was lowest for this category of hard negatives in Sugarcrepe [15]. We evaluate this paradigm on the swap-object split of Sugarcrepe and the VG-spatial split of VL-Checklist [72] against naive and finetuned CLIP. We choose this particular split of VL-Checklist because we find that all other splits (Objects, Attributes, or Relation-Action) do not introduce new functional words and are thus not applicable for testing dynamic FR updates.

All LLM queries were made to OpenAI's gpt4o-mini model, with a temperature of 0.7. At the time of evaluation, this was the most affordable model on the API at: \$ 0.150 / 1M input tokens and \$0.600 / 1M output tokens.

### F.1. LLM System Prompts

**System prompts:**

Given a list of sentences, reformat each sentence to the simplest phrases that would distinguish it from the other examples. Here are some formatting rules to follow:

1. If a sentence contains OBJECTS and ATTRIBUTES which belong to that object, the ATTRIBUTE must always come first. For example, given the sentence "A dog which is purple", reformat it to "A purple dog".

2. If a sentence contains NEGATION, the NEGATING TERM always comes before the OBJECT clause. For example, given "This image contains a chicken but a butterfly is absent", reformat it to "Chicken but no butterfly".

3. If a sentence contains PREPOSITIONs, try your best to make sure that the OBJECTS the PREPOSITION is describing are immediately before and after the PREPOSITION. For example, given "A bug which is flying much farther up from the bench", reformat it to "A flying bug above a bench".

4. Whenever possible, reformat VERBs to be ATTRIBUTES. For example, given "A dog dancing while his owner is jumping", reformat it to "Dancing dog and jumping owner".

5. If two sentences are very close to each other, reduce them down to the salient components. For example, given the sentences ["Butterflies in the clouds, a cat squatting looking up at it, and a man standing behind the cat watching it, on the grass with a tree.", "Butterflies in the clouds, a cat squatting looking up at it, and a man sitting behind the cat watching it, on the grass with a tree."], return: 0: "A man standing", 1: "A man sitting". (Of course, if there are other sentences in the list that are similar, you may want to keep more details so that the sentences are still distinguishable.)

Here are some more general examples. If given the following sentences: ["A desktop computer sitting on top of a gray oak table lights up the room" ,"A gray oak computer sitting on top of a desktop table lights up the room", "A kitchen has metal cabinets and black countertops with shiny lights on top.", "A kitchen has black cabinets and metal countertops with shiny lights on top."], return: 0: "Desktop computer top of gray oak table", 1: "Gray oak computer top of desktop table", 2: "Metal cabinets and black countertops", 3: "Black cabinets and metal countertops".

As a rule, be AS CONCISE AS POSSIBLE. If any information is repeated and unnecessary to keep in order to distinguish that text prompt from the others, discard it.

Now, reformat the following list of sentences and return the JSON output. Do not include anything other than the JSON array in your answer.

Table 7. Prompt template for simplifying natural language prompts.

**System prompts:**

You are given: A LOOKUP LIST of functional words (e.g., ["ABOVE", "BELOW", "INSIDE OF", "MANY", "SMALL", "NO"]). A list of SENTENCES to process. Definitions: Functional words include: (a) Prepositions (e.g., ABOVE, BELOW, INSIDE OF, ON, IN, NEAR) or their synonyms. (b) Size/shape terms (e.g., SMALL, BIG). (c) Numerical terms (e.g., ONE, TWO, THREE, etc.). If a number is greater than 5 (e.g., SEVEN, 100), replace it with "MANY". (d) Negatory terms (e.g., NO, WITHOUT). Non-functional words: Do not include verbs, adjectives, or any nouns unrelated to the functional categories above. Examples of non-functional words include "jumping", "sleeping", "cat", "man", etc. These should not be added to the LOOKUP LIST, even if they appear in the sentences. Even if there are two sentences that are very similar, do not try to distinguish them by adding these verbs, adjectives, or nouns to the LOOKUP LIST. Any form of a verb, including present participles, may not go in the LOOKUP LIST, no matter how frequently it appears in the sentences.

Rules: For each sentence: Identify any functional words or synonyms (including numbers). If a functional word or one of its synonyms (by meaning) appears in the sentence and is already in the LOOKUP LIST, replace it in the sentence with the LOOKUP LIST key, surrounded by angle brackets (e.g., "⟨ ABOVE⟩ "). If that functional word is not in the LOOKUP LIST (and it is truly functional by the above definition), add it to the LOOKUP LIST, then replace its appearance with that new all-caps key in angle brackets. Do not add duplicates to the LOOKUP LIST. Do not add verbs, adjectives, or any non-functional words to the LOOKUP LIST. Replace numbers greater than 5 with "MANY" (add "MANY" to the list if not already present). After processing all sentences, output exactly one JSON array containing two sub-arrays: The first sub-array: the UPDATED LOOKUP LIST (only functional words, no duplicates). The second sub-array: the FINAL TRANSFORMED SENTENCES (with functional words surrounded by ⟨ ⟩ ). Not every sentence needs functional words. Provide no additional commentary or text besides this JSON structure. Example of the required output format: [ [ "ABOVE", "INSIDE OF", "MANY" , "RIGHT OF"], [ "A bird ⟨ ABOVE⟩ a tree", "Fifteen dogs is ⟨ MANY⟩ dogs" , "A sitting chicken is ⟨ INSIDE OF⟩ a house"] ] Example of wrong output: [ [ "ABOVE", "INSIDE OF", "MANY" , "RIGHT OF", "SITTING", "DOGS", "HOUSE"], [ "A bird ⟨ ABOVE⟩ a tree", "Fifteen dogs is ⟨ MANY⟩ dogs" , "A sitting chicken is ⟨ INSIDE OF⟩ a house"] ]

Before you return the output, CHECK THAT THE LOOKUP LIST ONLY CONTAINS FUNCTIONAL WORDS.

Table 8. Prompt template for extracting functional words.

**Swapping Objects:**

Given an input sentence describing a scene, your task is to first locate two swappable noun phrases in the sentence, and then swap them to make a new sentence. The new sentence must meet the following three requirements: 1. The new sentence must be describing a different scene from the input sentence. 2. The new sentence must be fluent and grammatically correct. 3. The new sentence must make logical sense.

To complete the task, you should: 1. Answer the question of whether generating such a new sentence is possible using "Yes" or "No". 2. Output the swappable noun phrases. 3. Swap them to make a new sentence.

Please produce a **single JSON array** (no extra text or explanation) for each input sentence. If there are K input sentences, return a list with K JSON objects separated by commas. Each element in the array must be a JSON object with the following structure: { "possible": "⟨Yes or No⟩", "swappable-noun-phrases": ["⟨NP1⟩", "⟨NP2⟩"], "swapped-sentence": "⟨the swapped sentence⟩" }

Example JSON output for the original sentence: "A cat resting on a laptop next to a person." { "possible": "Yes", "swappable-noun-phrases": ["laptop", "person"], "swapped-sentence": "A cat resting on a person next to a laptop." }

Table 9. Prompt template for creating captions for swapped objects. This is a minorly changed version from [15].