

Towards Compositionality in Concept Learning

Adam Stein¹ Aaditya Naik¹ Yinjun Wu¹ Mayur Naik¹ Eric Wong¹

Abstract

Concept-based interpretability methods offer a lens into the internals of foundation models by decomposing their embeddings into high-level concepts. These concept representations are most useful when they are *compositional*, meaning that the individual concepts compose to explain the full sample. We show that existing unsupervised concept extraction methods find concepts which are not compositional. To automatically discover compositional concept representations, we identify two salient properties of such representations, and propose Compositional Concept Extraction (CCE) for finding concepts which obey these properties. We evaluate CCE on five different datasets over image and text data. Our evaluation shows that CCE finds more compositional concept representations than baselines and yields better accuracy on four downstream classification tasks.

1. Introduction

Foundation models continue to enable impressive performance gains across a variety of domains, tasks, and data modalities. However, their black-box nature severely limits the ability to debug, monitor, control, and trust them.

Concept-based explanations (Kim et al., 2018; Zhou et al., 2018) are a promising approach that seeks to explain a model’s behavior using individual concepts such as object attributes (e.g. *striped*) or linguistic sentiment (e.g. *happiness*). Decomposing a model’s learned representation can derive these concepts. For instance, a model’s embedding of a dog image may decompose into the sum of concept vectors representing its fur, snout, and tail.

Existing works based on methods such as PCA (Zou et al., 2023a) or KMeans (Ghorbani et al., 2019) extract such concept vectors reasonably well for basic concepts. For instance, Figure 1 shows images from the CUB (Wah et al., 2011) dataset containing concepts extracted by PCA from

¹Department of Computer and Information Science, University of Pennsylvania, Pennsylvania, USA. Correspondence to: Adam Stein <steinad@seas.upenn.edu>.

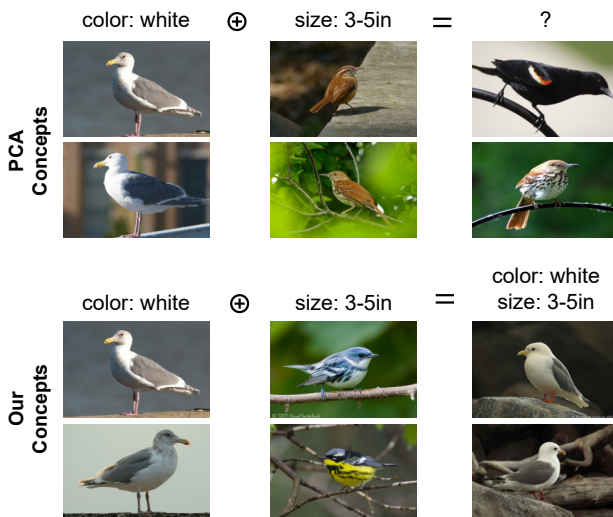


Figure 1. We illustrate the issue of concept compositionality with respect to concepts extracted from the embeddings of the CLIP model over the CUB dataset. Specifically, we visualize the concepts *white birds* and *small birds* learned by PCA (Zou et al., 2023a) and CCE along with their compositions. We show the top two images that best represent each concept. Ideally, composing the *white birds* and *small birds* concepts should result in a concept representing small white birds. This is not the case with the concepts learned by PCA. On the other hand, the concepts extracted by CCE are composable, as shown by the images of small white birds that best represent the resulting concept.

the CLIP (Radford et al., 2021) model. These techniques are able to correctly extract the representations of concepts like *white birds* and *small birds*, however, composing them by adding their representations does not yield the representation of the concept of *small white birds*.

The *compositionality* of concepts is vital for several use cases. First, model predictions can be explained by combining concepts (Abid et al., 2022). Compositional concepts also allow for editing fine-grained model behavior, like improving the truthfulness of an LLM without compromising other behaviors (Zou et al., 2023a). Models can also be trained to compose basic concepts for new tasks, e.g. using concepts for beak shapes, wing colors, and environments to classify bird species (Yuksekonul et al., 2023).

In this paper, we study the unsupervised extraction of compositional concepts. Existing work does not directly eval-

uate the compositionality of extracted concepts, but rather focuses on the individual concept representations. We therefore evaluate the compositionality of concepts extracted by existing unsupervised approaches.

For this purpose, we first validate the compositionality of ground-truth representations of concepts in controlled settings. We observe that concepts can be grouped into *attributes*, where each attribute consists of concepts over some common property, such as the color of objects or the shape of objects. Concepts from different attributes (e.g. `blue` and `cube`) can be composed, while those from the same attribute (e.g. `red` and `green`) cannot. We also observe that the concepts from different attributes are roughly orthogonal, while those from the same attribute are not. We prove in a generalized setting that these properties are crucial for the compositionality of concepts. Since existing approaches do not enforce these properties, they often extract non-composable concept representations.

To extract compositional concepts in an unsupervised manner, we propose Compositional Concept Extraction (CCE). Our key insight is to search for entire subspaces of concepts at once instead of individual concepts, allowing CCE to enforce the aforementioned properties of compositional concepts. We show that CCE recovers the representation of known compositional concepts better than existing approaches, can discover compositional concepts in existing image and text datasets, and the discovered concepts improve downstream classification accuracy.

We thus summarize the contributions of our paper:

- We study concept-based explanations of foundation models from the lens of compositionality—a property desirable for many use-cases. We observe that concept representations extracted by state-of-the-art methods fail to compose, and set out to remedy this problem.
- We validate that models can in fact represent concepts compositionally in embedding space. We identify two salient properties of compositional concept representations that existing methods fail to satisfy.
- We prove in a generalized setting that the identified properties are necessary for compositionality. We present a novel method called Compositional Concept Extraction (CCE) that guarantees to yield concept representations that satisfy these properties by construction.
- We demonstrate that CCE extracts more compositional concepts than baselines on vision and language datasets, and they improve downstream performance.

2. Concepts and Compositionality

Concept Representations. In machine learning, concepts are symbols that are assigned some human-interpretable meaning, often used to explain predictions made by models.

A concept extractor E extracts concepts from the intermediate representation of some pretrained model M over a dataset D . $E(M, D)$ thus yields a set of *concept vectors* representing the concepts $C = \{c_1, \dots, c_i\}$. Concept vectors are denoted as $R(c)$, where $R: \mathbb{C} \rightarrow \mathbb{R}^d$ is the concept representation function, \mathbb{C} is the set of all possible concepts, and \mathbb{R}^d is an embedding space in some dimension d . The set of extracted concepts C can be grouped into mutually exclusive *attributes* A_1, \dots, A_k each containing concepts about some common property such that $C = \bigcup_{i=1}^k A_i$.

To measure the presence (or degree of expression) of a concept in a sample’s embedding, we borrow the following definition of concept score from (Yeh et al., 2020).

Definition 2.1. (Concept Score) For a concept $c \in \mathbb{C}$ and concept representation function $R: \mathbb{C} \rightarrow \mathbb{R}^d$, a sample embedding $z \in \mathbb{R}^d$ has *concept score* $s(z, c) = S_{\cos}(z, R(c))$ where S_{\cos} is the *cosine similarity* function.

Existing work makes use of concept scores to quantify the presence of concepts on a per-sample basis. This has uses in several applications, such as creating concept bottleneck models where a sample embedding is converted to concept scores used for classification (Yuksekgonul et al., 2023), and sorting samples by a concept (Kim et al., 2018).

Compositionality. Following work on compositional representations (Andreas, 2019) and pretrained embeddings (Trager et al., 2023), we define the compositionality of concept representations.

Definition 2.2. (Compositional Concept Representations) For concepts $c_i, c_j \in \mathbb{C}$, the concept representation $R: \mathbb{C} \rightarrow \mathbb{R}^d$ is compositional if for some $w_{c_i}, w_{c_j} \in \mathbb{R}^+$,

$$R(c_i \cup c_j) = w_{c_i} R(c_i) + w_{c_j} R(c_j).$$

In other words, the representation of the composition of concepts corresponds to the weighted sum of the individual concept vectors in the embedding space.

Furthermore, concept scores for the concepts satisfying Definition 2.2 also behave compositionally, since each concept score quantifies the presence of that concept in a sample.

Lemma 2.3. For compositional concepts $c_i, c_j \in \mathbb{C}$, the concept score of their composition $c_k = c_i \cup c_j$ over a sample embedding $z \in \mathbb{R}^d$ is the composition of the concept scores of c_i and c_j , weighted by $w_{c_i}, w_{c_j} \in \mathbb{R}^+$:

$$s(z, c_k) = w_{c_i} s(z, c_i) + w_{c_j} s(z, c_j).$$

Since concept scores are used for several downstream tasks discussed above, this property about the compositionality of concept scores can simplify such tasks and improve the overall performance on them.

Besides finding compositional concepts, we also want to explain embeddings based on the concepts which compose

it. Prior work also performs a decomposition into a sum of concept representations (Zhou et al., 2018), but we modify the definition of such a decomposition so that a sample embedding is composed of only the concept representations that are truly present for the sample.

Definition 2.4. (Concept-based Decomposition) Consider a sample that is associated with a set of concepts $C \subseteq \mathbb{C}$, such that each attribute $A_i \subseteq C$ contains exactly one concept. A concept representation $R : \mathbb{C} \rightarrow \mathbb{R}^d$ decomposes that sample’s embedding $z_i \in \mathbb{R}^d$ if it can be expressed as the weighted sum of the sample’s associated concepts:

$$z_i = \sum_{c \in C} \lambda_{i,c} R(c), \text{ such that } \lambda_{i,j} > 0.$$

As an example, consider the CLEVR dataset (Johnson et al., 2017) consisting of images of objects of different shapes and colors. A concept extractor for a vision model may extract the set of concepts $C_{\text{CLEVR}} = \{\{\text{red}\}, \{\text{blue}\}, \{\text{cube}\}, \{\text{sphere}\}\}$. C_{CLEVR} can also be grouped into attributes $A_1 = \{\{\text{red}\}, \{\text{blue}\}\}$ and $A_2 = \{\{\text{cube}\}, \{\text{sphere}\}\}$ containing color and shape concepts respectively. As such, a composite concept like $\{\text{red}, \text{sphere}\}$ can be represented as the weighted sum of $R(\{\text{red}\})$ and $R(\{\text{sphere}\})$.

3. Evaluating Concept Compositionality

In this section, we validate the compositionality of ground-truth concept representations and evaluate the same for concepts extracted using existing approaches. We first discuss our controlled setting and show that concept representations from the CLIP model are compositional. We then evaluate the compositionality of concepts extracted by existing approaches. Finally, we outline the necessary properties of compositional concept representations.

3.1. Setup

In order to validate the compositionality of ground-truth concepts, we focus on concepts extracted from subsets of the CLEVR (Johnson et al., 2017), CUB (Wah et al., 2011), and Truth (Azaria & Mitchell, 2023) datasets, all of which have labelled attributes with compositional structure.

We follow a setup similar to (Lewis et al., 2022) for the synthetic CLEVR (Johnson et al., 2017) dataset and consider images with single objects labelled as one of three shapes (sphere, cube, or cylinder) and one of three colors (red, green, or blue). We also consider a subset of the CUB dataset consisting of bird images labelled as one of three colors and one of three sizes. Finally, we consider a subset of the Truth (Zou et al., 2023b) dataset consisting of facts relating to one of three topics and labelled true or false.

3.2. Ground-Truth Concept Compositionality

We evaluate the compositionality of ground-truth concept representations learned by the CLIP model over each labelled dataset. Since these representations are not provided, for each concept, we consider the mean of the model’s embeddings for samples belonging to that concept as a surrogate of its true representation (Zou et al., 2023a).

For example, for the CLEVR dataset, we extract the ground-truth representation of the `red` concept by calculating the mean of all sample embeddings of images with red objects. We similarly extract the ground-truth representations for the other two color concepts, the three shape concepts, and composite concepts like $\{\text{red}, \text{sphere}\}$, for a total of 15 concepts. We repeat this process for each dataset.

As stated in Lemma 2.3, the concept score for a composite of two concepts is the weighted sum of the concept scores of each concept. This implies that a linear model should be able to predict the concept score for a composed concept given the concept scores for each of the component concepts. We thus train a linear model to predict the presence or absence of a composed concept given its component concepts. We measure the average precision of the model for each composed concept, and report the mean average precision (MAP) score in Table 2a for each dataset. We see that in all cases, the ground truth (GT) concepts have high MAP (up to 0.971 for CLEVR) when predicting concept compositions from their components, meaning the ground-truth concept representations are reasonably compositional.

3.3. Compositionality Issues with Existing Methods

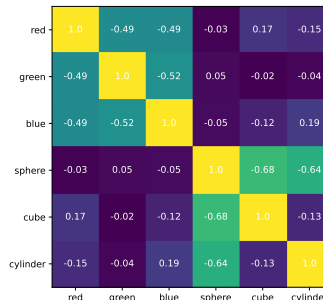
We next study the compositionality of concept representations discovered by existing unsupervised concept extraction methods. We train a linear model similar to the one described in Section 3.2, but with concepts extracted by baseline methods instead of the ground truths. From the MAP results in Table 2a we see that all the baselines have significantly lower compositionality than the ground-truth.

This is the case even for techniques that extract the concepts reasonably well, i.e. where the extracted concepts are able to discriminate between positive and negative samples of that concept. For each dataset and concept extraction method, we calculate the ROC-AUC score to measure the ability of the extracted concept to perform such a discrimination. We provide the full ROC-AUC results in Appendix E.1. In the case of NMF, despite this score averaging as high as 0.907 for the CLEVR dataset, the extracted concepts are not compositional. This implies that finding concept representations simply based on their ability to discriminate positive and negative samples of a concept does not mean that those representations will compose as expected.

We further demonstrate this point with a toy illustration in

Method	CLEVR	CUB-sub	Truth-sub
GT	0.995	0.808	0.728
PCA	0.964	0.602	0.467
ACE	0.683	0.670	0.592
Dictlearn	0.909	0.628	0.586
NMF	0.733	0.514	0.563
Concept Tf	0.499	0.502	0.378
Random	0.514	0.417	0.461
CCE	0.995	0.679	0.615

(a) MAP score of predicting concept compositions.



(b) Cosine similarities between CLEVR concepts.

Figure 2. Compositionality of ground-truth concepts compared with concepts extracted by existing approaches and CCE. Figure 2a shows that the ground-truth concepts (GT) are quite compositional, but existing methods are not. Figure 2b shows the cosine similarities between pairs of ground-truth concepts for the CLEVR dataset. The darker blue cells represent concepts that are orthogonal, while the lighter yellow ones represent non-orthogonal ones. We observe that concepts tend to be more orthogonal if they belong to different attributes.

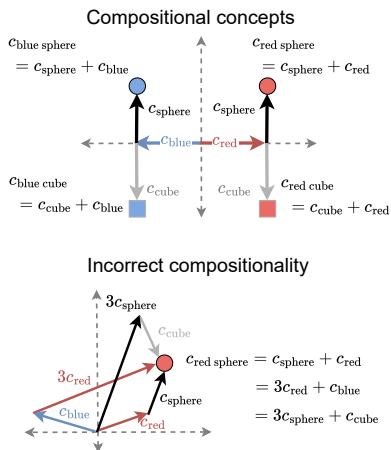


Figure 3. Illustration of concepts on a dataset of cubes and spheres that are either red or blue. The concepts on the top are compositional while those on the bottom are not. Even though the concepts on the bottom have perfectly represent the four samples, they still fail to compose properly. For instance, the composition of the red and blue concepts can form the {red, sphere} concept even though the blue concept is not present in a red sphere.

Figure 3. This figure depicts four perfectly composed concepts at the top, and four incorrectly composed concepts at the bottom, even though each concept is perfectly discriminative of the samples with the concept. Therefore, we must ensure that we explicitly extract compositional concepts.

3.4. Desired Properties of Compositional Concepts

To extract compositional concepts, we must first identify characteristics of such concepts. Since the ground-truth concepts were compositional, we investigate the salient characteristics of those concepts.

Consider the ground-truth concepts for the CLEVR dataset. In order to understand the relationship between different ground-truth concepts and their compositionality, we cen-

ter the sample embeddings and visualize cosine similarities between pairs of these concepts in Figure 2b. We observe that the ground-truth representations of color concepts are roughly orthogonal (cosine similarity near 0) to those of shape concepts. In contrast, the representations of concepts within the same attribute, such as the red and blue concepts, are non-orthogonal. Furthermore, the orthogonal concepts are also those that can compose to form new concepts, since they lie in different attributes. For instance, the red and sphere concepts are orthogonal, and can compose to form the {red, sphere} concept, while the red concept can't compose with the blue concept.

We visualize the same for the CUB-sub and Truth-sub datasets in Appendix C, and empirically observe the following trend over all three datasets: concept representations from different attributes are roughly orthogonal while those from the same attribute are non-orthogonal. Also, the orthogonal concepts tend to be compositional, while the non-orthogonal ones can't be composed.

Orthogonality is a generally helpful property for several use cases, such as disentangling concepts in embedding space (Chen et al., 2020). Some approaches therefore try to enforce orthogonality on the concepts being extracted. Table 1 summarizes existing unsupervised approaches for concept extraction and whether the method enforces any orthogonality constraints (Ortho.) between concepts of different attributes and allows for non-orthogonality between those of the same attribute (Corr.). We see that these approaches allow for only one of the two, but not both.

We now formally prove that the observed properties regarding concept compositionality hold in a generalized setting.

Theorem 3.1. For some dataset, let there be k attributes A_1, \dots, A_k and l concepts $c_{1,j}, \dots, c_{l,j}$ for each attribute A_j . Assuming that for each compositional concept $c = \{c_{1,j_1}, \dots, c_{k,j_k}\}$, its representation v_{j_1, \dots, j_k} , follows a

Table 1. Properties enforced by unsupervised concept extraction.

Method	Example	Ortho.	Corr.
PCA	RepE (Zou et al., 2023a)	✓	✗
KMeans	ACE (Ghorbani et al., 2019)	✗	✓
Dictionary-Learning	TransformerVis (Yun et al., 2021)	✗	✓
NMF	CRAFT (Fel et al., 2023)	✗	✓
Custom	Concept Tf (Rigotti et al., 2022)	✗	✓
Custom	CCE (Ours)	✓	✓

Algorithm 1 Compositional Concept Extraction

Input: embeddings Z , num. attr. M , concepts per attr. K
Initialize concepts $C = \{\}$
for $m = 1 \dots M$ **do**
 Initialize $P \in \mathbb{R}^{d \times k}$ such that $P^T P = I$.
 Initialize K concepts $V = \{v_1, \dots, v_K\}$.
 repeat
 $P = \text{LearnSubspace}(P, Z, V)$
 $V = \text{LearnConcepts}(ZP, K)$
 until Converged
 $C = C \cup V$
 $Z = Z - ZP$
end for
Return C

spherical normal distribution with zero mean and unit covariance, i.e. $v_{j1, \dots, jk} \sim N(\mathbf{0}, \mathbf{I}^d)$, the following statements are true with high probability for a large dimension d :

- For the base concepts belonging to the same attribute, there exists at least one pair of non-orthogonal concepts.
- For any pair of base concepts from two different attributes, they are orthogonal with high probability.

We show the proof in Appendix B. The takeaway from this result is that compositional concepts will be roughly orthogonal, while concepts of the same attribute may not be orthogonal. We leverage this to design an unsupervised concept extraction method which can find compositional concepts when they exist.

4. Compositional Concept Extraction (CCE)

To achieve this orthogonality property between concepts, we propose CCE, summarized in Algorithm 1 and visualized in Figure 4. As the outer loop of the algorithm suggests, once we find concepts for an attribute in a subspace P , we remove that subspace using orthogonal rejection and find concepts in a new subspace. This enforces the discovered subspaces to be orthogonal to each other, thus respecting the orthogonality property described in Section 3.

To discover concepts within each attribute, we employ a two-step process illustrated in Figure 4. The first step, shown on the left, is to project input data Z onto the subspace P_1

defined by the projection matrix $P_1 \in \mathbb{R}^{d \times s}$, resulting in projected data ZP_1 . In the next step, shown on the right, we identify concepts by performing spherical K-Means clustering within P_1 .

This clustering process is performed within a learned subspace. Therefore, we jointly learn the subspace P_1 and the clustering centroids within the subspace such that the data is highly clustered in that subspace. Specifically, we employ the Silhouette score (Rousseeuw, 1987) to quantify how well clustered the projected data ZP_1 is given some cluster assignment L determined by spherical K-means clustering. Since the Silhouette score is differentiable, we can search for an optimized P_1^* such that the Silhouette score is maximized:

$$P^* = \arg \max_{P, L} \text{Sil}(ZP, L).$$

We further observe that simply maximizing the above formula leads to overfitting issues since projecting the learned cluster centroids back to the original space may not necessarily correspond to cluster centroids in the original space. Therefore, we try to match the cluster centroids learned within the subspace and projected out to the original space to the centroids of the clusters in the original space. This is integrated into the above objective function as a regularization term, i.e.:

$$P^* = \arg \max_{P, L} \left(\text{Sil}(ZP, L) + \sum_k S_{\cos}(C_k P^T, \hat{C}_k) \right),$$

where C_k represents the clustering centroids in the subspace P_1 while $\hat{C}_k = \sum_i \mathbb{1}[L_i = k] Z_i$.

5. Experiments

5.1. Experimental Setup

Datasets and Models. We evaluate using five datasets across vision and language settings: CLEVR (Johnson et al., 2017) (vision), CUB (Wah et al., 2011) (vision), HAM10000 (Tschandl et al., 2018) (vision), Truth (Zou et al., 2023b) (language), and News (Mitchell, 1999) (language). We perform experiments on both *controlled* and *full* settings. In the controlled setting, we follow the same configuration as Section 3.1 for the CLEVR, CUB and Truth datasets. The *full setting* considers all samples from the CUB, Ham, Truth, and News datasets.

For the image datasets, we obtain sample representations from the CLIP model (Radford et al., 2021) while for the NLP dataset, this is achieved with Llama-2 13B Chat model (Touvron et al., 2023). We also perform ablation studies on the choices of different models in Appendix E.3.

Baseline Methods. Since the concept representations are learned by CCE in an unsupervised manner, we therefore primarily compare CCE against the following state-of-the-art unsupervised concept extraction methods, i.e., PCA (Zou

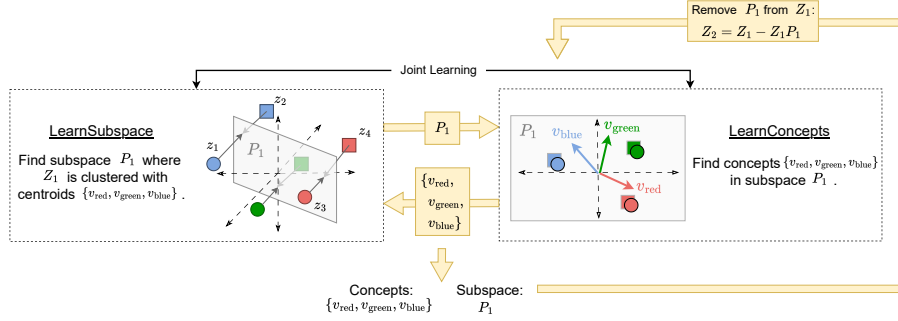


Figure 4. Finding color concepts in one iteration of CCE, which can be proceeded by finding other concepts, such as shapes.

et al., 2023a), NMF (Fel et al., 2023), ACE (KMeans) (Ghorbani et al., 2019), and Dictionary Learning (Bricken et al., 2023; Yun et al., 2021). In addition, we include a Random baseline where we randomly initialize concept vectors from a normal distribution of mean zero and variance one.

Recent studies like Concept Transformer (Rigotti et al., 2022) explore how to jointly learn concept representations and perform training of downstream classification tasks with learned concept representations. Hence, we treat Concept Transformer (Concept Tf) (Rigotti et al., 2022) as another baseline. Note that Concept Tf can optionally incorporate concept labels as additional supervisions, which are not considered in our experiments for fair comparison.

Experiment Design. We aim to answer these questions regarding the quality of the learned concept representations:

- RQ1** In the controlled setting with known compositional ground-truth concept representations, does CCE compose concepts more effectively than baselines?
- RQ2** In the full setting where the ground-truth concepts are typically unknown, can CCE successfully discover new and meaningful compositional concepts?
- RQ3** In both controlled and full settings, how can the learned compositional concept representations impact downstream performance?

To address **RQ1**, we evaluate the compositionality score (Andreas, 2019) on the concept representations extracted by CCE and the baselines, which is defined as follows:

Definition 5.1. (Compositionality Score) Given a dataset D consisting of embeddings $z \in \mathbb{R}^d$, their associated ground-truth concepts $C \subset \mathbb{C}$, and a concept representation function $R : \mathbb{C} \rightarrow \mathbb{R}^d$ obtained from a concept extractor, the compositionality score is the following:

$$\min_{\Lambda \geq 0} \frac{1}{|D|} \sum_{(z, C) \in D} \left\| z - \sum_{i=1}^{|C|} \Lambda_{z,i} R(C_i) \right\|$$

Intuitively speaking, for a sample embedding z , this metric quantifies how much z can be reconstructed by composing a list of concept representation $R(c_i)$'s that correspond to the i_{th} ground-truth concepts of z . Each $R(c_i)$ is weighted

Table 2. Compositionality Scores (lower is better).

	CLEVR	CUB-sub	Truth-sub
GT	3.05	0.475	3.74
PCA	3.64	0.480	4.06
ACE	3.59	0.527	3.75
Dictlearn	3.43	0.502	3.75
NMF	3.60	0.542	3.99
Concept Tf	4.89	0.546	4.93
Random	4.93	0.546	4.35
CCE	3.06	0.474	3.74

by a coefficient $\Lambda_{z,i}$, which is determined by optimizing the above formula with respect to all $\Lambda_{z,i}$.

In addition, for each ground-truth concept, we also report the cosine similarity between the learned concept representation $R(c_i)$ and the corresponding ground-truth representation.

To study **RQ2** for the full setting, we primarily perform qualitative studies to identify whether CCE is capable of discovering reasonable compositional concepts. Specifically, for each learned concept representation, assign a name to the concept by inspecting the ten images with the top concept score. Then for each pair of the learned concepts, we first identify those samples with the highest concept scores. Then, we sum the two concept representations, and find the samples with largest concept score for this aggregated representation. By investigating these examples, we visually examine whether the composition is reasonable or not.

Lastly, we answer **RQ3** by evaluating the downstream classification performance with the learned concept representations. Specifically, we follow Yuksekogonul et al. (2023) to learn a linear classifier by predicting class labels with the concept scores of a sample. We further report the performance of training a linear classifier on sample embeddings without involving any concepts, denoted by ‘‘No concept’’.

5.2. Experimental Results

Compositionality in Controlled Settings.

We first evaluate the compositionality scores on the CLEVR,

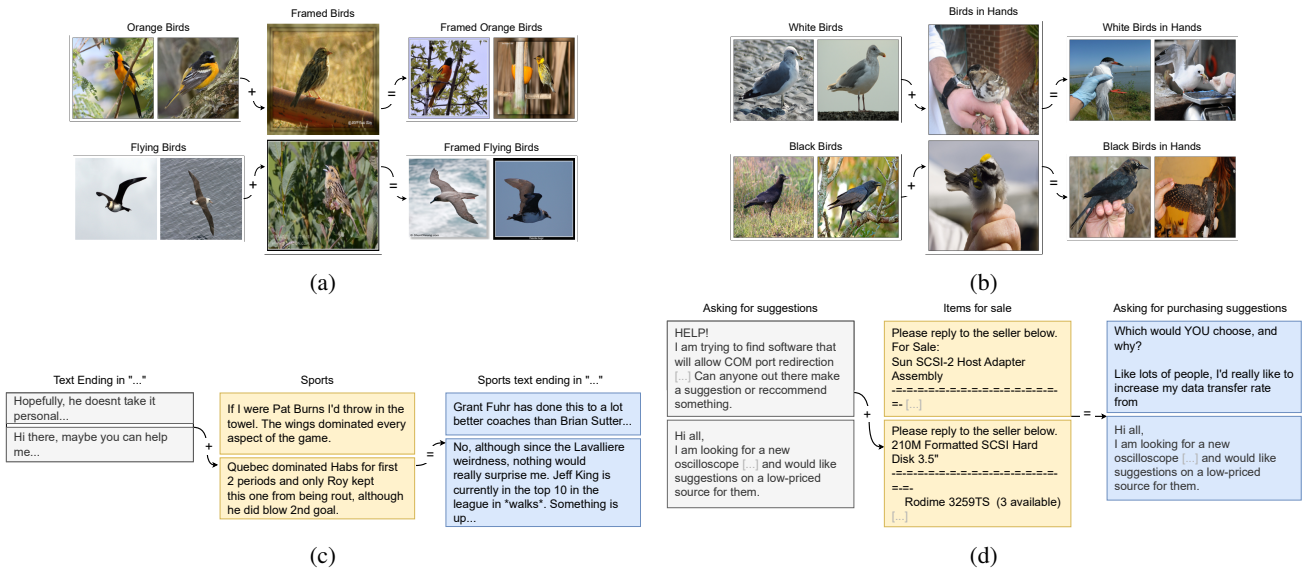


Figure 5. Examples of compositional concepts identified by CCE. Figures 5a and 5b are from the CUB dataset while Figures 5c and 5d are from the News dataset. These figures suggest that CCE can not only discover new meaningful concepts outside the ground-truth concepts, such as the Birds in Hands concept in Figure 5a, but also compose these concepts correctly, e.g. White Birds + Birds in Hands = White Birds in Hands.

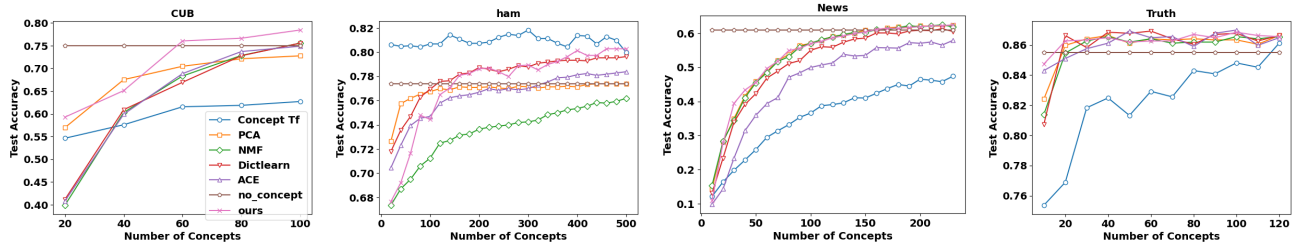


Figure 6. Downstream classification accuracy on the full setting.

Table 3. The average cosine similarity between individual learned concept representations and the ground truth (higher is better).

	CLEVR	CUB-sub	Truth-sub
PCA	0.608	0.503	0.459
ACE	0.760	0.724	0.716
Dictlearn	0.775	0.653	0.633
NMF	0.738	0.116	0.651
Concept Tf	0.058	0.061	0.012
Random	0.044	0.076	0.019
CCE	0.988	0.749	0.876

CUB-sub, and Truth-sub datasets and report them in Table 2. In all cases, CCE obtains the best score compared to the baselines, indicating the advantage of CCE in discovering compositional concepts. Moreover, CCE’s scores are comparable to those of the ground-truth concept representations. This shows that the concepts learned by CCE almost align with the ground-truth concept representations.

This is further supported by the results in Table 3. This table summarizes the cosine similarities between the ground-truth

concept representations and the ones learned by the baselines and CCE. Again, the concepts learned by CCE are the closest to the ground truths. Note that some baselines like Dictlearn also produce highly accurate concept representations. However, as Table 2 shows, their compositions fail to be consistent with the ground truths.

Compositionality in Real Data Settings. To address RQ2, we perform some qualitative studies on compositional concepts discovered by CCE on the CUB and News dataset, which are visualized in Figure 5. As shown in this figure, CCE is capable of identifying reasonable concepts, such as White Birds, Framed Birds and Text Ending in ‘‘...’’. Some of these concepts are even beyond the ground-truth concept labels that are provided by the dataset itself. For example, CCE identifies the ‘‘Birds in Hands’’ concept which is not labeled in the CUB dataset. But its top activated samples are images with a bird in someone’s hand (see Figure 5a). Furthermore, the composition of those learned concepts is also representative of the properties of each concept. For example, in Figure 5c, the composition of the concept Text Ending in ‘‘...’’

and `Sports` represents sentences about “sports” ending in “...”.

Downstream Performance Analysis. For `RQ3`, we studied the impact of the extracted compositional concepts on downstream performance across all datasets in the full setting. Throughout the experiments, we observe that the total number of concepts is a crucial factor in determining the performance. Therefore, we also vary this number and report the performance numbers accordingly for all datasets and methods in Figure 6. As this figure suggests, across all the datasets, despite the poor performance with a small number of concepts, CCE gradually gains performance with an increasing number of concepts, eventually outperforming all the unsupervised baseline methods.

Also, it is worth noting that CCE outperforms Concept Tf most times and is on par with it in the worst case (see the experimental results on the ham dataset with 500 concepts). This thus indicates the performance advantage of CCE even in the absence of supervision from downstream tasks.

Furthermore, CCE discovers concept representations by performing a series of linear transformations on top of the sample embeddings. But by comparing against “No concept” where sample embeddings are directly used for downstream tasks, CCE can even outperform it by a large margin on CUB and Ham dataset. This implies that the concept representations extracted by CCE might be more relevant to the downstream classification tasks than the raw embeddings.

6. Related Work

Concept-based Interpretability. Concept-based interpretability encompasses the building of models using human-interpretable concepts (Koh et al., 2020; Espinosa Zarlenga et al., 2022; Yuksekgonul et al., 2023) and extracting such concepts post-hoc from models (Kim et al., 2018; Zhou et al., 2018). In either case, how do we choose which concepts to use? Some existing work specifies concepts using human supervision to select and provide their labels (Kim et al., 2018), large-scale concept annotation datasets (Bau et al., 2017), general knowledge bases (Yuksekgonul et al., 2023), and large language models (Yang et al., 2023). Another line of work uses regularization (Wong et al., 2021) or other inductive biases (Rigotti et al., 2022) to learn concepts during standard supervised training of a model. Finally, another line of work leverages unsupervised methods to automatically discover concepts (Ghorbani et al., 2019; Fel et al., 2023; Yun et al., 2021; Bricken et al., 2023). This paper focuses finding concepts from a pretrained model without supervision.

Compositionality in Foundation Models. Since the observation of compositional word vectors by Mikolov et al. (2013) there has been interest in finding and utilizing compo-

sitional behavior of deep learning models. Compositionality has been used to uncover and mitigate bias in word embeddings (Bolukbasi et al., 2016), edit classifier behavior (Santurkar et al., 2021), and recently to monitor and control the behavior of foundational language (Todd et al., 2023; Zou et al., 2023a) and vision models (Wang et al., 2023; Kwon et al., 2023). In this work, we identify compositionality of concepts in datasets for pretrained vision and language models, and how to automatically discover such concepts.

Compositional and Disentangled Representations. In representation learning, there is considerable effort to encourage *disentangled* representations (Bengio et al., 2013; Higgins et al., 2016; Wang et al., 2022). While disentanglement concerns how to distinguish separate concepts in embedding space, compositionality concerns what happens when separate concepts get combined. Existing work has shown that disentanglement and compositionality do not have to be correlated (Xu et al., 2022). Unlike representation learning, we start with a pretrained model and try to uncover the compositional concepts it learned.

7. Limitations

We study the case where concepts compose compositionally, but concepts may also be non-compositional. For instance, the concepts of `hot` and `dog` do not compose to form the meaning of `hot dog` (Zhai, 1997). In addition, we supposed a flat concept structure, which does not distinguish between “(small blue) car” and “small (blue car)”. We leave the study of such non-compositional and hierarchical concepts to future work.

Another limitation of unsupervised concept extraction is that discovered concept vectors are not associated with any name. We assign names to the concept through manual inspection of samples with a high concept score, but this can require significant effort with large numbers of concepts.

8. Conclusion

In this paper, we studied concept-based explanations of foundation models from the lens of compositionality. We validated that the ground-truth concepts extracted from these models are compositional while the existing unsupervised concept extraction methods usually fail to guarantee compositionality. To address this issue, we first identified two salient properties for compositional concept representations and designed a novel concept extraction method called CCE that respects these properties by design. Through extensive experiments across vision and language datasets, we demonstrated that CCE not only learns compositional concepts but also enhances downstream performance.

9. Broader Impacts

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Abid, A., Yuksekgonul, M., and Zou, J. Meaningfully debugging model mistakes using conceptual counterfactual explanations. In *International Conference on Machine Learning*, pp. 66–88. PMLR, 2022.
- Andreas, J. Measuring compositionality in representation learning. In *International Conference on Learning Representations*, 2019.
- Azaria, A. and Mitchell, T. The internal state of an llm knows when its lying. *arXiv preprint arXiv:2304.13734*, 2023.
- Bau, D., Zhou, B., Khosla, A., Oliva, A., and Torralba, A. Network dissection: Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 6541–6549, 2017.
- Bengio, Y., Courville, A., and Vincent, P. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8): 1798–1828, 2013.
- Bolukbasi, T., Chang, K.-W., Zou, J. Y., Saligrama, V., and Kalai, A. T. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29, 2016.
- Bricken, T., Templeton, A., Batson, J., Chen, B., Jermyn, A., Conerly, T., Turner, N., Anil, C., Denison, C., Askell, A., Lasenby, R., Wu, Y., Kravec, S., Schiefer, N., Maxwell, T., Joseph, N., Hatfield-Dodds, Z., Tamkin, A., Nguyen, K., McLean, B., Burke, J. E., Hume, T., Carter, S., Henighan, T., and Olah, C. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023. <https://transformer-circuits.pub/2023/monosemantic-features/index.html>.
- Chen, Z., Bei, Y., and Rudin, C. Concept whitening for interpretable image recognition. *Nature Machine Intelligence*, 2(12):772–782, 2020.
- Espinosa Zarlenga, M., Barbiero, P., Ciravegna, G., Marra, G., Giannini, F., Diligenti, M., Shams, Z., Precioso, F., Melacci, S., Weller, A., et al. Concept embedding models: Beyond the accuracy-explainability trade-off. *Advances in Neural Information Processing Systems*, 35:21400–21413, 2022.
- Fel, T., Picard, A., Bethune, L., Boissin, T., Vigouroux, D., Colin, J., Cadène, R., and Serre, T. Craft: Concept recursive activation factorization for explainability. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2711–2721, 2023.
- Ghorbani, A., Wexler, J., Zou, J. Y., and Kim, B. Towards automatic concept-based explanations. *Advances in neural information processing systems*, 32, 2019.
- Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., and Lerchner, A. beta-vae: Learning basic visual concepts with a constrained variational framework. In *International conference on learning representations*, 2016.
- Johnson, J., Hariharan, B., Van Der Maaten, L., Fei-Fei, L., Lawrence Zitnick, C., and Girshick, R. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2901–2910, 2017.
- Kim, B., Wattenberg, M., Gilmer, J., Cai, C., Wexler, J., Viegas, F., et al. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *International conference on machine learning*, pp. 2668–2677. PMLR, 2018.
- Koh, P. W., Nguyen, T., Tang, Y. S., Musmann, S., Pierson, E., Kim, B., and Liang, P. Concept bottleneck models. In *International conference on machine learning*, pp. 5338–5348. PMLR, 2020.
- Kwon, M., Jeong, J., and Uh, Y. Diffusion models already have a semantic latent space. In *The Eleventh International Conference on Learning Representations*, 2023.
- Lewis, M., Nayak, N. V., Yu, P., Yu, Q., Merullo, J., Bach, S. H., and Pavlick, E. Does clip bind concepts? probing compositionality in large image models. *arXiv preprint arXiv:2212.10537*, 2022.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26, 2013.
- Mitchell, T. Twenty Newsgroups. UCI Machine Learning Repository, 1999. DOI: <https://doi.org/10.24432/C5C323>.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.

- Rigotti, M., Mikšović, C., Giurgiu, I., Gschwind, T., and Scotton, P. Attention-based interpretability with concept transformers. In *International conference on learning representations*, 2022.
- Rousseeuw, P. J. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics*, 20:53–65, 1987.
- Santurkar, S., Tsipras, D., Elango, M., Bau, D., Torralba, A., and Madry, A. Editing a classifier by rewriting its prediction rules. *Advances in Neural Information Processing Systems*, 34:23359–23373, 2021.
- Todd, E., Li, M. L., Sharma, A. S., Mueller, A., Wallace, B. C., and Bau, D. Function vectors in large language models. *arXiv preprint arXiv:2310.15213*, 2023.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Trager, M., Perera, P., Zancato, L., Achille, A., Bhatia, P., and Soatto, S. Linear spaces of meanings: compositional structures in vision-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 15395–15404, 2023.
- Tschandl, P., Rosendahl, C., and Kittler, H. The ham10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. *Scientific data*, 5(1):1–9, 2018.
- Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. The caltech-ucsd birds-200-2011 dataset. 2011.
- Wang, X., Chen, H., Tang, S., Wu, Z., and Zhu, W. Disentangled representation learning. *arXiv preprint arXiv:2211.11695*, 2022.
- Wang, Z., Gui, L., Negrea, J., and Veitch, V. Concept algebra for (score-based) text-controlled generative models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Wegner, S.-A. Lecture notes on high-dimensional data. *arXiv preprint arXiv:2101.05841*, 2021.
- Wong, E., Santurkar, S., and Madry, A. Leveraging sparse linear layers for debuggable deep networks. In *International Conference on Machine Learning*, pp. 11205–11216. PMLR, 2021.
- Xu, Z., Niethammer, M., and Raffel, C. A. Compositional generalization in unsupervised compositional representation learning: A study on disentanglement and emergent language. *Advances in Neural Information Processing Systems*, 35:25074–25087, 2022.
- Yang, Y., Panagopoulou, A., Zhou, S., Jin, D., Callison-Burch, C., and Yatskar, M. Language in a bottle: Language model guided concept bottlenecks for interpretable image classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19187–19197, 2023.
- Yeh, C.-K., Kim, B., Arik, S., Li, C.-L., Pfister, T., and Ravikumar, P. On completeness-aware concept-based explanations in deep neural networks. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 20554–20565. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/ecb287ff763c169694f682af52c1f309-Paper.pdf.
- Yuksekgonul, M., Wang, M., and Zou, J. Post-hoc concept bottleneck models. In *The Eleventh International Conference on Learning Representations*, 2023.
- Yun, Z., Chen, Y., Olshausen, B., and Lecun, Y. Transformer visualization via dictionary learning: contextualized embedding as a linear superposition of transformer factors. In *Proceedings of Deep Learning Inside Out (DeeLIO): The 2nd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, pp. 1–10, 2021.
- Zhai, C. Exploiting context to identify lexical atoms—a statistical view of linguistic context. *arXiv preprint cmp-lg/9701001*, 1997.
- Zhou, B., Sun, Y., Bau, D., and Torralba, A. Interpretable basis decomposition for visual explanation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 119–134, 2018.
- Zou, A., Phan, L., Chen, S., Campbell, J., Guo, P., Ren, R., Pan, A., Yin, X., Mazeika, M., Dombrowski, A.-K., et al. Representation engineering: A top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*, 2023a.
- Zou, A., Phan, L., Chen, S., Campbell, J., Guo, P., Ren, R., Pan, A., Yin, X., Mazeika, M., Dombrowski, A.-K., et al. Representation engineering: A top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*, 2023b.

A. Proof of Lemma 2.3

Proof. Let $z \in \mathbb{R}^d$ be a sample embedding, $R : \mathbb{C} \rightarrow \mathbb{R}^d$ be a compositional concept representation function, and $c_i, c_j \in \mathbb{C}$ be two compositional concepts which compose as $c_k = c_i \cup c_j$. From Definition 5.1, the concept scores for c_i and c_j are the following:

$$\begin{aligned} s(z, c_i) &= S_{\cos}(z, R(c_i)) \\ s(z, c_j) &= S_{\cos}(z, R(c_j)). \end{aligned}$$

The concept score for the composition c_k can then be written as:

$$\begin{aligned} s(z, c_k) &= s(z, c_i \cup c_j) \\ &= S_{\cos}(z, R(c_i \cup c_j)) \\ &= S_{\cos}(z, w_{c_i} R(c_i) + w_{c_j} R(c_j)) && \text{(since } R \text{ is compositional)} \\ &= \frac{z \cdot (w_{c_i} R(c_i) + w_{c_j} R(c_j))}{\|z\| \|w_{c_i} R(c_i) + w_{c_j} R(c_j)\|} && \text{(definition of cosine similarity)} \\ &= \frac{z \cdot w_{c_i} R(c_i)}{\|z\| \|R(c_k)\|} + \frac{z \cdot w_{c_j} R(c_j)}{\|z\| \|R(c_k)\|} \\ &= \frac{(w_{c_i} \|R(c_i)\|) z \cdot R(c_i)}{\|R(c_k)\| \|z\| \|R(c_i)\|} + \frac{(w_{c_j} \|R(c_j)\|) z \cdot R(c_j)}{\|R(c_k)\| \|z\| \|R(c_j)\|} \\ &= \frac{w_{c_i} \|R(c_i)\|}{\|R(c_k)\|} S_{\cos}(z, R(c_i)) + \frac{w_{c_j} \|R(c_j)\|}{\|R(c_k)\|} S_{\cos}(z, R(c_j)) && \text{(definition of cosine similarity)} \end{aligned}$$

□

B. Proof of Theorem 3.1

Lemma B.1 (curse of dimensionality). (*Wegner, 2021*) For a pair of vectors \mathbf{x} and \mathbf{y} randomly sampled from $N(0, \mathbf{I}^d)$, \mathbf{x} and \mathbf{y} are orthogonal with high probability for large enough d . Mathematically speaking, for a fixed small constant, ϵ , the following inequality holds:

$$\mathbb{P} \left[\left| \left\langle \frac{\mathbf{x}}{\|\mathbf{x}\|}, \frac{\mathbf{y}}{\|\mathbf{y}\|} \right\rangle \right| \leq \epsilon \right] \geq 1 - \frac{M_1}{\sqrt{d}\epsilon} - \frac{M_2}{\sqrt{d}},$$

where $M_1 = 2$ and $M_2 = 7$

Lemma B.2 (Gaussian Annulus Theorem). (*Wegner, 2021*) For a vector v randomly sampled from $N(0, \mathbf{I}^d)$, $\|v\|$ is approaching \sqrt{d} with high probability for large enough d . Mathematically speaking, the following inequality holds:

$$\mathbb{P} \left[\left| \|\mathbf{x}\| - \sqrt{d} \right| \leq \epsilon \right] \geq 2 \exp(-M_3 \epsilon^2),$$

in which $M_3 = \frac{1}{16}$

Based on the above two lemmas, for any two randomly sampled vectors \mathbf{x} and \mathbf{y} from $N(0, \mathbf{I}^d)$, the following equality holds with high probability:

$$\langle \mathbf{x}, \mathbf{y} \rangle = o(d) \tag{1}$$

Proof. First, we can derive the concept representation for each base concept $c_{i,t}$ (denoted by $\mu_{i,t}$) as follows:

$$\mu_{i,t} = \frac{1}{l^{k-1}} \sum_{j_1} \sum_{j_2} \cdots \sum_{j_{i-1}} \sum_{j_{i+1}} \cdots \sum_{j_k} v_{j_1, j_2, j_3, \dots, j_{i-1}, t, j_{i+1}, \dots, j_k}. \tag{2}$$

Since we also want to perform centering operations over the entire dataset, then this suggests that we need to leverage the mean of all concepts, i.e.,:

$$\mu = \frac{1}{l^k} \sum_{j_1} \sum_{j_2} \cdots \sum_{j_k} v_{j_1, j_2, j_3, \dots, j_{i-1}, j_i, j_{i+1}, \dots, j_k}. \quad (3)$$

Then after the centering operation, $\mu_{i,t}$ is transformed into:

$$\mu'_{i,t} = \frac{\mu_{i,t} - \mu}{\sigma'}.$$

In the formula above, we use σ' to represent the standard deviation vector calculated over the entire dataset.

Then let us fix i and sum up all $\mu'_{i,t}$ over all t , which yields:

$$\begin{aligned} \sum_{t=1}^l \mu'_{i,t} &= \sum_{t=1}^l \frac{\mu_{i,t} - \mu}{\sigma'} \\ &= \sum_{t=1}^l \frac{\mu_{i,t}}{\sigma'} - \frac{l \cdot \mu}{\sigma'} \end{aligned} \quad (4)$$

Then by integrating Equation equation 2 and Equation equation 3 into the above formula, we can get:

$$\begin{aligned} \sum_{t=1}^l \mu'_{i,t} &= \frac{1}{\sigma'} \left[\sum_{t=1}^l \frac{1}{l^{k-1}} \sum_{j_1} \sum_{j_2} \cdots \sum_{j_{i-1}} \sum_{j_{i+1}} \cdots \sum_{j_k} v_{j_1, j_2, j_3, \dots, j_{i-1}, t, j_{i+1}, \dots, j_k} \right. \\ &\quad \left. - \frac{1}{l^{k-1}} \sum_{j_1} \sum_{j_2} \cdots \sum_{j_{i-1}} \sum_{j_{i+1}} \cdots \sum_{j_k} v_{j_1, j_2, j_3, \dots, j_{i-1}, t, j_{i+1}, \dots, j_k} \right] \\ &= \frac{1}{\sigma'} \left[\frac{1}{l^{k-1}} \sum_{j_1} \sum_{j_2} \cdots \sum_{j_{i-1}} \sum_{t=1}^l \sum_{j_{i+1}} \cdots \sum_{j_k} v_{j_1, j_2, j_3, \dots, j_{i-1}, t, j_{i+1}, \dots, j_k} \right. \\ &\quad \left. - \frac{1}{l^{k-1}} \sum_{j_1} \sum_{j_2} \cdots \sum_{j_{i-1}} \sum_{j_{i+1}} \cdots \sum_{j_k} v_{j_1, j_2, j_3, \dots, j_{i-1}, t, j_{i+1}, \dots, j_k} \right] = 0. \end{aligned}$$

This thus suggests that for a fixed i , all $\mu'_{i,t}$ are correlated. If all pairs of $(\mu'_{i,t_1}, \mu'_{i,t_2}), t_1, t_2 = 1, 2, \dots, l$ are orthogonal, then we can obtain the following formula:

$$\langle \mu'_{i,t_1}, \sum_{t=1}^l \mu'_{i,t} \rangle = \langle \mu'_{i,t_1}, \mu'_{i,t_1} \rangle = 0,$$

thus indicating that all $\mu_{i,t}$ is 0. This is thus contradictory to our assumption that each $\mu_{i,t}$ is a non-zero vector. Therefore, there should exist at least one pair of $(\mu'_{i,t_1}, \mu'_{i,t_2})$ which are not orthogonal.

We next prove that arbitrary pairs of concept representations from two different attributes are orthogonal with high probability. To demonstrate this, we calculate the dot product between μ'_{i_1, t_1} and μ'_{i_2, t_2} which represents two concepts from attribute i_1 and i_2 respectively:

$$\begin{aligned} \langle \mu'_{i_1, t_1}, \mu'_{i_2, t_2} \rangle &= \left\langle \frac{\mu_{i_1, t_1} - \mu}{\sigma'}, \frac{\mu_{i_2, t_2} - \mu}{\sigma'} \right\rangle \\ &= \frac{1}{\sigma'^2} \frac{1}{l^k} \left\langle l \sum_{j_1} \sum_{j_2} \cdots \sum_{j_{i_1-1}} \sum_{j_{i_1+1}} \cdots \sum_{j_k} v_{j_1, j_2, j_3, \dots, j_{i_1-1}, t_1, j_{i_1+1}, \dots, j_k} - \sum_{j_1} \sum_{j_2} \cdots \sum_{j_k} v_{j_1, j_2, j_3, \dots, j_k}, \right. \\ &\quad \left. l \sum_{j_1} \sum_{j_2} \cdots \sum_{j_{i_2-1}} \sum_{j_{i_2+1}} \cdots \sum_{j_k} v_{j_1, j_2, j_3, \dots, j_{i_2-1}, t_2, j_{i_2+1}, \dots, j_k} - \sum_{j_1} \sum_{j_2} \cdots \sum_{j_k} v_{j_1, j_2, j_3, \dots, j_k} \right\rangle \end{aligned}$$

According to Equation equation 1, for arbitrary pairs of v_{j_1, j_2, \dots, j_k} and $v_{j'_1, j'_2, \dots, j'_k}$, as long as their indexes are not exactly equivalent, their dot product is $o(d)$. Therefore, through some linear algebraic operations, the above formula could be reformulated as follows:

$$\begin{aligned} \langle \mu'_{i_1, t_1}, \mu'_{i_2, t_2} \rangle &= \frac{1}{\sigma' l^2} \left(l^2 \sum_{j_1} \sum_{j_2} \cdots \sum_{j_{i_1-1}} \sum_{j_{i_1+1}} \cdots \sum_{j_{i_2-1}} \sum_{j_{i_2+1}} \cdots \sum_{j_k} \|v_{j_1, j_2, j_3, \dots, j_{i_1-1}, t_1, j_{i_1+1}, \dots, j_{i_2-1}, t_2, j_{i_2+1}, \dots, j_k}\|_2^2 \right. \\ &\quad - l \sum_{j_1} \sum_{j_2} \cdots \sum_{j_{i_1-1}} \sum_{j_{i_1+1}} \cdots \sum_{j_k} \|v_{j_1, j_2, j_3, \dots, j_{i_1-1}, t_1, j_{i_1+1}, \dots, j_k}\|_2^2 \\ &\quad - l \sum_{j_1} \sum_{j_2} \cdots \sum_{j_{i_2-1}} \sum_{j_{i_2+1}} \cdots \sum_{j_k} \|v_{j_1, j_2, j_3, \dots, j_{i_2-1}, t_2, j_{i_2+1}, \dots, j_k}\|_2^2 \\ &\quad \left. + \sum_{j_1} \sum_{j_2} \cdots \sum_{j_k} \|v_{j_1, j_2, j_3, \dots, j_k}\|_2^2 \right) + o(d) \end{aligned}$$

Plus, according to Lemma B.2, for each vector x randomly sampled from $N(\mathbf{0}, \mathbf{I}^d)$, its norm is bounded by $[\sqrt{d} - \epsilon, \sqrt{d} + \epsilon]$ with high probability, which applies to each $v_{j_1, j_2, j_3, \dots, j_k}$. Then the above formula could be further bounded by:

$$\begin{aligned} \langle \mu'_{i_1, t_1}, \mu'_{i_2, t_2} \rangle &\leq \frac{1}{\sigma' l} \left(l^2 \sum_{j_1} \sum_{j_2} \cdots \sum_{j_{i_1-1}} \sum_{j_{i_1+1}} \cdots \sum_{j_{i_2-1}} \sum_{j_{i_2+1}} \cdots \sum_{j_k} (\sqrt{d} + \epsilon)^2 \right. \\ &\quad - l \sum_{j_1} \sum_{j_2} \cdots \sum_{j_{i_1-1}} \sum_{j_{i_1+1}} \cdots \sum_{j_k} (\sqrt{d} - \epsilon)^2 \\ &\quad - l \sum_{j_1} \sum_{j_2} \cdots \sum_{j_{i_2-1}} \sum_{j_{i_2+1}} \cdots \sum_{j_k} (\sqrt{d} - \epsilon)^2 \\ &\quad \left. + \sum_{j_1} \sum_{j_2} \cdots \sum_{j_k} (\sqrt{d} + \epsilon)^2 \right) + o(d) \\ &= \frac{1}{\sigma' l} \left[l^k (\sqrt{d} + \epsilon)^2 - l^k (\sqrt{d} - \epsilon)^2 - l^k (\sqrt{d} - \epsilon)^2 + l^k (\sqrt{d} + \epsilon)^2 \right] + o(d) \\ &= \frac{1}{\sigma'} \left[8l^{k-1} \sqrt{d} \epsilon \right] + o(d), \end{aligned}$$

Similarly, we can prove that

$$\langle \mu'_{i_1, t_1}, \mu'_{i_2, t_2} \rangle \geq -\frac{1}{\sigma'} \left[8l^k \sqrt{d} \epsilon \right] + o(d)$$

Therefore, we can conclude that

$$\langle \mu'_{i_1, t_1}, \mu'_{i_2, t_2} \rangle = o(d) \quad (5)$$

In addition, we can compute the norm of μ'_{i_1, t_1} and follow the same derivation as above by leveraging Equation equation 1, which results in:

$$\begin{aligned} \|\mu'_{i_1, t_1}\|_2^2 &= l \sum_{j_1} \sum_{j_2} \cdots \sum_{j_{i_1-1}} \sum_{j_{i_1+1}} \cdots \sum_{j_k} v_{j_1, j_2, j_3, \dots, j_{i_1-1}, t_1, j_{i_1+1}, \dots, j_k} - \sum_{j_1} \sum_{j_2} \cdots \sum_{j_k} v_{j_1, j_2, j_3, \dots, j_k} \|_2^2 \\ &= l^2 \sum_{j_1} \sum_{j_2} \cdots \sum_{j_{i_1-1}} \sum_{j_{i_1+1}} \cdots \sum_{j_k} \|v_{j_1, j_2, j_3, \dots, j_{i_1-1}, t_1, j_{i_1+1}, \dots, j_k}\|_2^2 + \sum_{j_1} \sum_{j_2} \cdots \sum_{j_k} \|v_{j_1, j_2, j_3, \dots, j_k}\|_2^2 + o(d). \end{aligned}$$

This formula could then be lower bounded by:

$$\|\mu'_{i_1, t_1}\|_2^2 \geq 2l^k (d - 2\sqrt{d}\epsilon + \epsilon^2) + o(d) = 2l^k d + o(d) \quad (6)$$

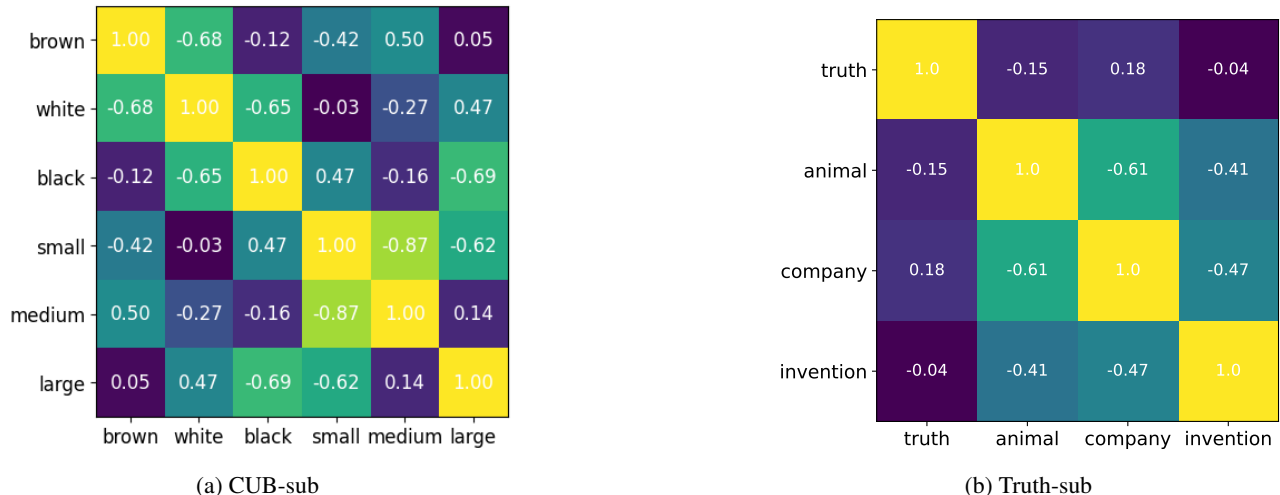


Figure 7. Compositionality of Ground-Truth Concepts for the CUB-sub and Truth-sub datasets.

This leverages the fact that each $\|v_{j_1, j_2, j_3, \dots, j_k}\|$ is bounded by $[\sqrt{d} - \epsilon, \sqrt{d} + \epsilon]$ with high probability. The above formula also holds for $\|\mu'_{i_2, t_2}\|_2^2$. As a consequence, the cosine similarity between μ'_{i_1, t_1} and μ'_{i_2, t_2} is bounded by:

$$\text{cosine}(\mu'_{i_1, t_1}, \mu'_{i_2, t_2}) = \frac{\langle \mu'_{i_1, t_1}, \mu'_{i_2, t_2} \rangle}{\|\mu'_{i_1, t_1}\| \cdot \|\mu'_{i_2, t_2}\|} \leq \frac{o(d)}{2l^k d + o(d)},$$

which thus approaches zero as d increases. □

C. Compositionality of Ground-Truth Concepts

D. Qualitative Examples

E. Additional quantitative examples

E.1. ROC-AUC Scores between Concept Representations and Ground-Truth

The maximum ROC-AUC between the concept score and the true label for the ground-truth concepts is presented in Table 4 for CLEVR, Table 5 for CUB-sub, and Table 6 for Truth-sub.

E.2. The analysis of the cosine similarity score between learned concept representations and ground-truth

We further break down the results reported in Table 3 average cosine similarity between the learned concept representation and the ground-truth concept representations.

E.3. Ablation studies on other pretrained models

Recall that in the experiment section, we primarily focus on discovering concepts from pretrained CLIP model. In this section, we study with different choices of pretrained models, can we obtain similar results as that in Section 5?

To answer this question, we leverage vision transformer (ViT), another widely used pretrained vision model, to repeat the experiments on CLEVR dataset. The results are summarized in Table 10-11. The results from these tables maintain the same trends as the one shown in Section 5.

Table 4. Max AUC score CLEVR v/s GT

Concepts	ours	ace	ace-svm	pca	dictlearn	seminmf
red	1.000	0.765	0.728	0.985	0.757	0.793
green	1.000	0.771	0.711	0.996	0.797	0.818
blue	1.000	0.753	0.745	0.972	0.782	0.836
sphere	1.000	1.000	0.736	1.000	1.000	1.000
cube	1.000	0.998	0.742	0.971	0.994	0.999
cylinder	1.000	0.998	0.831	0.977	0.992	0.998
(red and sphere) object	0.987	0.993	0.911	0.950	0.978	0.983
(red and cube) object	0.923	0.999	1.000	0.965	0.983	0.999
(red and cylinder) object	0.899	0.940	0.932	0.964	0.998	0.943
(green and sphere) object	0.858	0.991	0.870	0.863	0.980	0.986
(green and cube) object	0.878	1.000	1.000	0.877	0.951	1.000
(green and cylinder) object	0.936	0.916	0.960	0.969	1.000	0.994
(blue and sphere) object	0.952	0.996	1.000	0.834	0.940	0.997
(blue and cube) object	0.878	1.000	1.000	0.973	0.842	0.978
(blue and cylinder) object	0.923	0.992	1.000	0.990	0.995	0.995

Table 5. ROC AUC of baseline methods on recovering the labeled concepts.

Method	Brown	White	Black	Small	Medium	Large
GT	0.984	0.999	0.998	1.000	0.923	0.847
PCA	0.881	0.985	0.931	0.997	0.886	0.677
ACE	0.895	0.785	0.677	0.726	0.584	0.678
DictLearn	0.849	0.645	0.650	0.702	0.519	0.551
NMF	0.086	0.164	0.099	0.116	0.066	0.168
Concept Tf.	0.923	0.837	0.887	0.926	0.754	0.736
Random	0.867	0.933	0.855	0.888	0.849	0.723
CCE	0.894	0.834	0.710	0.743	0.656	0.661

Table 6. ROC AUC of baseline methods on recovering the labeled concepts.

Method	Truth	Animal	Company	Invention
GT	0.91	1.00	1.00	1.00
PCA	0.829	0.917	0.832	0.863
ace	0.777	0.999	0.941	0.795
dictlearn	0.353	0.734	0.627	0.539
nmf	0.759	0.708	0.629	0.521
Ours	0.91	1.00	0.96	0.78

Table 7. Max AUC score CLEVR v/s GT ViT

Concepts	ours	ace	ace-svm	pca	dictlearn	seminmf
red	1.000	0.688	0.735	0.945	0.710	0.712
green	1.000	0.692	0.711	0.922	0.716	0.680
blue	1.000	0.692	0.642	0.995	0.704	0.629
sphere	1.000	1.000	0.610	1.000	1.000	1.000
cube	1.000	1.000	0.735	0.970	0.999	1.000
cylinder	1.000	1.000	0.695	1.000	1.000	1.000
(red and sphere) object	0.972	0.998	1.000	0.980	0.997	0.991
(red and cube) object	0.884	0.986	0.720	0.881	0.992	0.967
(red and cylinder) object	0.933	0.935	0.837	0.962	0.998	1.000
(green and sphere) object	0.904	1.000	1.000	0.923	0.998	0.985
(green and cube) object	0.913	0.992	0.731	0.886	0.920	0.937
(green and cylinder) object	0.895	0.905	0.660	0.866	0.988	0.939
(blue and sphere) object	0.939	0.960	0.844	0.970	0.954	0.949
(blue and cube) object	0.825	0.847	0.770	0.905	0.838	0.851
(blue and cylinder) object	0.854	0.899	0.766	0.842	0.913	0.875

Table 8. Cosine Similarity results on CLEVR

Concepts	ours	ACE	PCA	Dictlearn	NMF	Concept Tf	Random
red	0.982	0.598	0.723	0.624	0.576	0.078	0.056
green	0.999	0.564	0.825	0.652	0.678	0.060	0.030
blue	0.973	0.659	0.075	0.699	0.786	0.057	0.059
sphere	1.000	0.912	0.995	0.916	0.813	0.043	0.030
cube	1.000	0.847	0.712	0.874	0.787	0.082	0.054
cylinder	1.000	0.925	0.150	0.805	0.747	0.028	0.036
(red and sphere) object	0.844	0.983	0.839	0.953	0.953		
(red and cube) object	0.847	1.000	0.769	0.971	0.927		
(red and cylinder) object	0.801	0.827	0.530	0.985	0.832		
(green and sphere) object	0.843	0.950	0.851	0.908	0.865		
(green and cube) object	0.790	0.999	0.597	0.889	0.950		
(green and cylinder) object	0.747	0.836	0.634	0.978	0.877		
(blue and sphere) object	0.857	0.981	0.842	0.849	0.884		
(blue and cube) object	0.773	0.999	0.318	0.669	0.844		
(blue and cylinder) object	0.850	1.000	0.155	0.980	0.959		

Table 9. Cosine similarity of baseline methods for recovering the labeled concepts.

Method	Truth	Animal	Company	Invention
pca	0.469	0.072	0.730	0.564
ace	0.640	0.854	0.728	0.641
dictlearn	0.322	0.866	0.724	0.619
nmf	0.428	0.862	0.770	0.542
Concept tf				
Random				
CCE	0.514	0.996	0.995	0.998

Table 10. Compositionality results on CLEVR (ViT)

Method	Compositionality
ours	3.688
ace	4.264
pca	4.367
dictlearn	3.757
nmf	4.202
Concept tf	
Random	

Table 11. Cosine Similarity results on CLEVR ViT

Concepts	ours	ACE	PCA	Dictlearn	NMF	Concept Tf	Random
red	0.997	0.410	0.686	0.436	0.441		
green	0.997	0.465	0.663	0.462	0.353		
blue	0.989	0.412	0.515	0.368	0.365		
sphere	1.000	0.883	0.982	0.863	0.882		
cube	1.000	0.969	0.077	0.974	0.899		
cylinder	1.000	0.934	0.923	0.921	0.951		
(red and sphere) object	0.872	1.000	0.861	0.996	0.970		
(red and cube) object	0.907	0.867	0.348	0.848	0.893		
(red and cylinder) object	0.849	0.867	0.794	0.994	0.956		
(green and sphere) object	0.881	1.000	0.853	0.994	0.936		
(green and cube) object	0.877	0.936	0.289	0.895	0.874		
(green and cylinder) object	0.881	0.822	0.793	0.840	0.882		
(blue and sphere) object	0.840	0.958	0.833	0.943	0.814		
(blue and cube) object	0.894	0.891	0.212	0.871	0.817		
(blue and cylinder) object	0.885	0.853	0.826	0.921	0.880		

Table 12. Cosine similarity of baseline methods for recovering the labeled concepts.

Method	Truth	Animal	Company	Invention
PCA	0.367	0.139	0.688	0.583
ace	0.244	0.956	0.733	0.642
dictlearn	0.760	0.988	0.917	0.879
nmf	0.824	0.898	0.931	0.725
Ours	0.90	0.94	0.85	0.64