Abstract

Interest in understanding and factorizing learned embedding spaces through conceptual explanations is steadily growing. When no human concept labels are available, concept discovery methods search trained embedding spaces for interpretable concepts like *object shape* or *color* that can be used to provide post-hoc explanations for decisions. Unlike previous work, we argue that concept discovery should be *identifiable*, meaning that a number of known concepts can be provably recovered to guarantee reliability of the explanations. As a starting point, we explicitly make the connection between concept discovery and classical methods like Principal Component Analysis and Independent Component Analysis by showing that they can recover independent concepts with non-Gaussian distributions. For dependent concepts, we propose two novel approaches that exploit functional compositionality properties of image-generating processes. Our provably identifiable concept discovery methods substantially outperform competitors on a battery of experiments including hundreds of trained models and dependent concepts, where they exhibit up to 29 % better alignment with the ground truth. Our results provide a rigorous foundation for reliable concept discovery without human labels.

1 INTRODUCTION

Modern computer vision systems represent and reason about images in embedding spaces. These are either constructed implicitly in higher-level layers of large models or explicitly through generative models such as Variational Autoencoders (Kingma and Welling, 2013) or Diffusion Models (Song and Ermon, 2019; Ho et al., 2020). To unveil why an image is considered similar to a certain class, interest in

understanding these embeddings is increasing. Conceptual explanations (Crabbé and van der Schaar, 2022; Muttenthaler et al., 2022; Akula et al., 2020; Kazhdan et al., 2020; Yeh et al., 2019; Kim et al., 2018) are a popular explainable AI (XAI) technique for this purpose. They scrutinize a given encoder by decomposing its embedding space into interpretable concepts post-hoc, i.e., after training. Subsequently, these concepts form the bases of popular post-hoc explanations such as TCAV (Kim et al., 2018) or allow high-level interventions (Koh et al., 2020). Fig. 1 outlines a real-world example. A misclassification made by a pretrained model shipped with the pytorch library (Paszke et al., 2017) is to be explained. In the given example, the conceptual explanation allows identification of a spurious correlation, that the model has picked up: Most jack-o-lanterns are found in combination with dark backgrounds, which causes it to mistake the traffic light at night for a jack-o-lantern.

Constructing such explanations is non-trivial. The key ingredient to all conceptual explanation techniques is a set of interpretable concepts, which is notoriously hard to specify (Leemann et al., 2022). It is frequently defined through human annotations (Crabbé and van der Schaar, 2022; Koh et al., 2020; Kim et al., 2018) on individual samples of the dataset that can be prohibitively expensive (Kazhdan et al., 2021). Furthermore, it is usually unknown which concepts will be leveraged by a machine learning model without a model at hand. Therefore, we consider fully unsupervised concept discovery (Ghorbani et al., 2019; Yeh et al., 2019), where the concepts are automatically discovered in the data. Concepts are frequently modeled as directions in a given embedding space (Ghorbani et al., 2019; Kim et al., 2018; Yeh et al., 2019), that have to be discovered without supervision. These embedding spaces can be highly distorted, making it hard to correctly separate the influences of individual concepts. However, this is essential to make the right inferences in practice (see Fig. 1d).

While many methods have been empirically shown to work well, a rigorous theoretical analysis of the conditions under which concept discovery is possible is still lacking in



(a) Misclassification: A this incident happened.

(b) Conceptual Explanation:

(c) **Inspection:** A closer inspection model makes an incorrect Concept contributions are com- of samples from the predicted class prediction. A user is inter- puted that explain the prediction. reveals that most images in this class ested to understand why In this example, the concept "dark- have a dark background; a spurious ness" is relevant for the outcome. correlation picked up by the model.

(d) Entangled Conceptual Explanation: It is essential to correctly split up the contribution of individual concepts to allow for valid inferences.

Figure 1: Schematical use-case of conceptual explanations: A misclassification of an image classifier is explained. The example is based on a real explanation for a ResNet50 model. Details and the original explanation are provided in App. C.8.

previous works. We propose to consider concept discovery methods that are *identifiable*. This means when a known number of ground truth components generated the data, the concept discovery method provably yields concepts that correspond to the individual ground truth components and can correctly represent an input in the concept space. This is a crucial requirement: If a method is even incapable of recovering known components, there is no indication for why it should be reliable in practice. In this work, we are the first to investigate identifiability results in the context of post-hoc concept discovery.

First, we find that identifiability results from Principal Component Analysis (PCA) and Independent Component Analysis (ICA) literature (Jolliffe, 2002; Comon, 1994; Hyvärinen et al., 2001) can be transferred to the conceptual explanation setup. We establish that they cover the case of independent ground truth components with non-Gaussian distributions. This is insufficient for two reasons: (1) In practice, concepts such as height and weight (Träuble et al., 2021) or wing and head colors of birds often follow complex dependency patterns. (2) Popular generative models (Kingma and Welling, 2013; Song and Ermon, 2019) frequently work with an embedding space with a Gaussian distribution.

As a second contribution, we seek to fill this void by providing an identifiable concept discovery approach that can handle dependent and Gaussian ground truth components. We can show that this is possible through taking the nature of the image-generating process into consideration. Specifically, we propose utilizing visual compositionality properties. These are based on the observation that tiny changes in the components frequently affect input images in orthogonal or even disjoint ways. These properties of image-generating processes also leave a "trace" in the encoders learned from a set of data samples. This insightful finding permits to construct two novel post-hoc concept discovery methods based on the *disjoint* or *independent mechanisms* criterion. We prove strong identifiability guarantees for recovering components, even if they are dependent.

In summary, our work advances current literature in multiple ways: (1) We present first identifiability results for post-hoc conceptual explanations. We find that results from ICA can be transferred under independent ground truth components. (2) For the more intricate setting of dependent components, we derive the *disjoint mechanism analysis (DMA)* and the less constrained independent mechanism analysis (IMA) criterion. We prove that they recover even dependent original components up to permutation and scale. (3) We construct DMA and IMA-based concept discovery algorithms for encoder embedding spaces with the same theoretical identifiability guarantees. (4) We test them (i) on embeddings of several autoencoder models learned from correlated data, (ii) with multiple and strong correlations, (iii) on discriminative encoders, and (iv) on the real-world CUB-200-2011 dataset (Wah et al., 2011). Our approaches maintain superior performance amidst increasingly severe challenges.

2 **RELATED WORK**

Works on the analysis and interpretation of embedding spaces touch a variety of subfields of machine learning.

Concept discovery for explainable AI. Conceptual explanations (Koh et al., 2020; Kim et al., 2018; Ghorbani et al., 2019; Yeh et al., 2019; Akula et al., 2020; Chen et al., 2020b) have gained popularity within the XAI community. They aim to explain a trained machine learning model post-hoc in terms of human-friendly, high-level concept directions (Kim et al., 2018). These concepts are found via supervised (Koh et al., 2020; Kim and Mnih, 2018; Kazhdan et al., 2020) or unsupervised approaches (Yeh et al., 2019; Akula et al., 2020; Ren et al., 2022), such as clustering of latent representations (Ghorbani et al., 2019). However, their results are not always meaningful (Leemann et al., 2022; Yeh et al., 2019). Therefore, we suggest approaches with identifiability guarantees. We provide initial identifiability results and a novel approach which can be used for unsupervised concept discovery under correlated components.

Independent Component Analysis (ICA). Independent Component Analysis (Comon, 1994; Hyvärinen and Pajunen, 1999; Hyvärinen et al., 2001) or blind source separation (BSS) view g(z) as a mixture to undo and rely on traces that the distributions over z leave in the mixture. In this work, we show that an identifiability result from ICA can be transferred to the conceptual explanation setup, but recovery is only possible under independent underlying components of which all but one are non-Gaussian. This result is not applicable to naturally correlated processes, which is why we design a novel method for this case.

Disentanglement Learning. Concurrently, literature in on disentanglement learning is concerned with finding a datagenerating mechanism g(z) and a latent representation zfor a dataset, such that each of the original components (also known as factors of variation) is mapped to one (controllable) unit direction in z (Bengio et al., 2013). An alternative definition relies on group theory (Higgins et al., 2017) where certain group operations (symmetries) should be reflected in the learned representation (Painter et al., 2020; Yang et al., 2021). Most works in the domain enhance VAEs (Kingma and Welling, 2013) with additional loss terms (Higgins et al., 2017; Burgess et al., 2018; Kim and Mnih, 2018; Chen et al., 2018). Despite recent progress it is not always possible to construct disentangled embedding spaces from scratch: Locatello et al. (2019) have shown that the problem is inherently unidentifiable without additional assumptions. A more recent work by Träuble et al. (2021) shows that even if just two components of a dataset are correlated, current disentanglement learning methods fail. In this work, we focus on post-hoc explanations of embedding spaces of given models, which are usually entangled.

Identifiability results. A strain of works have considered identifiability in disentanglement learning. It has been previously shown that unsupervised disentanglement, without further condition, is impossible (Hyvärinen and Pajunen, 1999; Locatello et al., 2019). Hence, recent works aim to understand the conditions sufficient for identifiability. One strain of work relies on additional supervision, i.e., access to an additional observed variable (Hyvärinen et al., 2019; Khemakhem et al., 2020) or to tuples of observations that differ in only a limited number of components (Locatello et al., 2020). Gresele et al. (2021) and Zheng et al. (2022) proved identifiable disentanglement under independently distributed components and introduce a functional condition on the data generator. We also consider functional properties but our setting is different as (1) we have access to a trained encoder only and (2) not even partial annotations or relations are available.

3 ANALYSIS

In this section, we formalize post-hoc concept discovery to provide an identifiability perspective. We find that Independent Component Analysis (ICA) and Principal Component Analysis (PCA) only guarantee identifiability when the ground-truth components are stochastically independent. We then study the intricate case of dependent components and propose using *disjoint* and *independent mechanisms analysis* (DMA / IMA) along with identifiability results. All proofs are provided in the supplementary.

3.1 PROBLEM FORMALIZATION

In post-hoc concept discovery, we are given a trained encoder $\boldsymbol{f}: \mathcal{X} \to \mathcal{E}$ with embeddings $\boldsymbol{e} = \boldsymbol{f}(\boldsymbol{x}) \in \mathcal{E} \subset \mathbb{R}^{K}$ of each image $x \in \mathcal{X}$. We do not impose any restriction on how f was obtained; it can be a the feature extractor part of a large classification model or a feature representation learned through autoencoding, constrastive learning (Chen et al., 2020a) or related techniques. Interpretability literature seeks to understand the embedding space by factorizing it into concepts. Based on the observations that directions in the embedding space often correspond to meaningful features (Szegedy et al., 2013; Bau et al., 2017; Alain and Bengio, 2016; Bisazza and Tump, 2018), these concepts are frequently defined as direction vectors m_i (Kim et al., 2018; Ghorbani et al., 2019; Yeh et al., 2019). Hence, the combined output of a concept discovery algorithm is a matrix $\boldsymbol{M} = [\boldsymbol{m}_1, \dots, \boldsymbol{m}_K]^{ op} \in \mathbb{R}^{K imes K}$ where each row contains a concept direction.

We seek a theoretical guarantee on when these discovered concept directions align with ground truth components that generated the data. To this end, we formalize the data-generating process as shown in Fig. 2: There are K ground-truth components with scores $z_k, k = 1 \dots K$, summarized $z \in \mathcal{Z} \subset \mathbb{R}^{K}$, that define an image. The term components always refers to the ground truth as opposed to the concepts, which denote the discovered directions. A data-generating process $g:\mathcal{Z}
ightarrow \mathcal{X}$ generates images $\boldsymbol{x} = \boldsymbol{g}(\boldsymbol{z}) \in \mathcal{X} \subset \mathbb{R}^L, L \gg K.$ A powerful algorithm should be able to recover the original components. That is, there should be a one-to-one mapping between entries of Me and the entries in z, up to the arbitrary scale and order of the entries. We say that a concept discovery algorithm identifies the true components if it is guaranteed to output directions M that satisfy Me = Mf(g(z)) = PSz $\forall z \in \mathcal{Z}$, where $P \in \mathbb{R}^{K \times K}$ is a permutation matrix that has one 1 per row and column and is 0 otherwise, and $\boldsymbol{S} \in \mathbb{R}^{K \times K}$ is an invertible diagonal scaling matrix.

To make the problem solvable in the first place, concept directions must exist in the embedding space of the given encoder, requiring e = Dz, where $D \in \mathbb{R}^{K \times K}$ is full-rank. Depending on the scope of the conceptual explanation desired, it can be sufficient for the components to exist in a local region of the embedding space if the concept discovery algorithm is only applied around a region around a certain point of interest. This only changes the meaning



Figure 2: Overview over the concept discovery setup. We consider a process where data samples x are generated from possibly correlated ground truth components z, e.g., a wingspan or beak length of a bird, by an unknown process g (left). The high-dimensional data is mapped to the to the embedding space of a given model f (center). A suitable post-hoc concept discovery yields concept vectors m_i that correspond to the original components (right).

of \mathcal{E}, \mathcal{X} , and \mathcal{Z} , but is formally equivalent.

3.2 IDENTIFIABILITY VIA INDEPENDENCE

Initially, we turn towards classical component analysis methods. We find that they require non-correlation or stronger stochastic independence of the ground truth components.

Principal Component Analysis (PCA) (Jolliffe, 2002) uses eigenvector decompositions to find orthogonal directions Mthat result in uncorrelated components Me. This means that PCA is only capable of identifying the original components if the ground truth components z were uncorrelated and exist as orthogonal directions in our embedding space. In our setup and notation, this leads to the following result:

Theorem 3.1 (PCA identifiability) Let $z_k, k = 1, ..., K$, be uncorrelated random variables with non-zero and unequal variances. Let e = Dz, where $D \in \mathbb{R}^{K \times K}$ is an orthonormal matrix. If an orthonormal post-hoc transformation $M \in \mathbb{R}^{K \times K}$ results in mutually uncorrelated components $(z'_1, ..., z'_K) = z' = Me$, then Me = PSz, where $P \in \mathbb{R}^{K \times K}$ is a permutation and $S \in \mathbb{R}^{K \times K}$ is a diagonal matrix where $|s_{ii}| = 1$ for $i \in 1, ..., K$.¹

All proofs in this work are deferred to App. B. It is arguably a strong condition that the ground truth directions are encoded orthogonally in the embedding space. Independent Component Analysis (ICA) overcomes this limitation Comon (1994) and allows for arbitrary directions, but requires stochastically independent components instead of the weaker non-correlation.

Theorem 3.2 (ICA identifiability) Let $z_k, k = 1, ..., K$, be independent random variables with non-zero variances where at most one component is Gaussian. Let e = Dz, where $D \in \mathbb{R}^{K \times K}$ has full rank. If a post-hoc transformation $M \in \mathbb{R}^{N \times N}$ results in mutually independent compo-

Dependency	Marginal Dist.	Transform	Criterion
uncorr.	uneq. variances	orthogonal	non-correlation (PCA)
independent	non-Gaussian	invertible	independence (ICA)
arbitrary	arbitrary	invertible	disj. mechanisms (DMA)
arbitrary	arbitrary	invertible	indep. mechanisms (IMA)

Table 1: PCA and ICA provably identify concepts via their distributions. DMA and IMA utilize functional properties.

nents
$$(z'_1, \ldots, z'_K) = z' = Me$$
, then $Me = PSz$, where $P \in \mathbb{R}^{K \times K}$ is a perm. and $S \in \mathbb{R}^{K \times K}$ is a diag. matrix.

This result shows that stochastic independence of the ground truth components leaves a strong trace in the embeddings that can be leveraged. Algorithms like fastICA (Hyvärinen and Oja, 1997) can find the concept directions M by searching for independence (Comon, 1994). We conclude that ICA is suited for post-hoc concept discovery under independent components.

In summary, we have transferred two results from the component analysis literature to the setup of post-hoc conceptual explanations. However these results do not allow to recover components that are correlated or follow a Gaussian distribution. This limits their applicability in practice where concepts often appear pairwise (e.g., darkness and jack-olanterns, cf. Fig. 1). We will bridge this gap in the remainder of this paper by introducing two new identifiable discovery methods based on functional properties of the generation process that we term *disjoint* and *independent* mechanisms. A summary of identifiability results is provided in Table 1.

3.3 IDENTIFIABILITY VIA DISJOINT MECHANISMS

Instead of placing independence assumptions on z, we propose a concept discovery algorithm that makes use of natural properties of the generative process g. In particular, generative processes in vision are often compositional (Ommer and Buhmann, 2007): Different groups of pixels in an image,

¹To simplify notation, P and S mean *any* permutation and scale matrices. They do not have to be equal between the theorems.

like a bird's wings, legs, and head, are each controlled by different components. Effects of tiny changes in components are visible the Jacobian J_g , where each row points to the pixels affected. Thus, a compositional process will follow the *disjoint mechanisms* principle.

Definition 3.1 (Disjoint mechanisms) g is said to generate x from its components z via disjoint mechanisms if the Jacobian $J_g(z) \in \mathbb{R}^{L \times K}$ exists and is a block matrix $\forall z \in \mathcal{Z}$. That is, the columns of $J_g(z)$ are nonzero at disjoint rows, i.e. $|J_g(z)|^\top |J_g(z)| = S(z)$, where $S \in \mathbb{R}^{K \times K}$ is a diagonal matrix that may be different for each z and $|\cdot|$ takes the element-wise absolute value.

Note that this definition does not globally constrain the location of affected pixels such that components may be alter, different, but disjoint pixels in each image. In real concept discovery, we do not have access to the generative process g but can only access the encoder f. However, an encoder corresponding to g will not be arbitrary and its Jacobian $J_f \in \mathbb{R}^{K \times L}$ have a distinct form in practice: First, to maintain the component information the composition $f \circ g$ will be of the form $\boldsymbol{f}(\boldsymbol{g}(\boldsymbol{z})) = \boldsymbol{D}\boldsymbol{z},$ with a yet unknown matrix $D \in \mathbb{R}^{K \times K}$. Furthermore, we expect encoders to be rather lazy, meaning they only perform the changes to invert the data generation process but are almost invariant to input deviations not due to changes in the components. Technically, the changes effected by the components form the linear span $(J_q(z))$, whereas entirely external changes are given in its orthogonal complement span $(J_g(z))^{\perp}$. Thus, for $\mathbf{v} \in \operatorname{span}(\boldsymbol{J}_{\boldsymbol{g}}(\boldsymbol{z}))^{\perp} \subset \mathbb{R}^{\hat{L}}$ the encoder should not react to these, i.e., $J_f(z)\mathbf{v} = \mathbf{0} \Leftrightarrow \mathbf{v} \in \ker(J_f(z)).$

Definition 3.2 (Faithful encoder) f is a faithful encoder for the generative process g if the ground truth components remain recoverable, i.e., f(g(z)) = Dz, for some $D \in \mathbb{R}^{K \times K}$ with full rank. Furthermore, f is lazy and invariant to changes in x which cannot be explained by the ground truth components, requiring $J_f(g(z))$ and $J_g(z)$ to exist and span $(J_g(z))^{\perp} \subseteq \ker(J_f(z)), \forall z \in \mathbb{Z}$.

Having defined what realistic encoders look like, we find, there is distinct property which can be leveraged to discover the directions in M among faithful encoders: It is necessary to find a decoder Mf whose Jacobian MJ_f will have disjoint rows. Intuitively, this requires searching for components whose gradients affect disjoint image regions.

Theorem 3.3 (Identifiability under DMA) Let g have disjoint mechanisms and f be a faithful encoder to g. If a full-rank post-hoc transformation $M \in \mathbb{R}^{K \times K}$ results in disjoint rows in the Jacobian $MJ_f(g(z))$, i.e., $|MJ_f(g(z))||MJ_f(g(z))|^{\top}$ is invertible and diagonal for some $z \in Z$, then Me = PSz, where $P \in \mathbb{R}^{K \times K}$ is a permutation and $S \in \mathbb{R}^{K \times K}$ is a scaling matrix. This theorem does not impose any restrictions on the distribution z, making it applicable to realistic concept discovery scenarios through leveraging the nature of the generative process. The proof of this algorithm in App. B.5 also yields an analytical solution. We will use it to verify conditions in a controlled experiment in Sec. 4.1. We have thus identified the *DMA criterion* that allows to discover the component directions: The rows of MJ_f need to point to disjoint image regions. We can formulate this as a loss function and optimize for M via off-the-shelf gradient descent:

$$\mathcal{L}(\boldsymbol{M}) = \mathbb{E}_{\boldsymbol{x}} \| \operatorname{arn} \left[\boldsymbol{M} \boldsymbol{J}_{f}(\boldsymbol{x}) \right] \operatorname{arn} \left[\boldsymbol{M} \boldsymbol{J}_{f}(\boldsymbol{x}) \right]^{\top} - \boldsymbol{I} \|_{F}^{2}.$$
(1)

The expectation is taken over a collection of real data samples x. The arn-operator (<u>absoute values</u>, <u>row normalization</u>) takes the element-wise absolute value and subsequently normalizes the rows. This does not constrain the norms of the Jacobian's rows but only enforces disjointness.

3.4 CONCEPT DISCOVERY VIA INDEPENDENT MECHANISMS

We can perform an analogous derivation for a class of generating processes that is more general. Grounded by causal principles instead of compositionality, the independent mechanisms property has been argued to define a class of natural generators (Gresele et al., 2021).

Definition 3.3 (Independent mechanisms (IMA)) g is said to generate x from its components z via independent mechanisms if the Jacobian $J_g(z)$ of g exists and its columns (one per component) are orthogonal $\forall z \in \mathbb{Z}$, i.e., $J_g^{\top}(z)J_g(z) = S(z)$, where $S \in \mathbb{R}^{K \times K}$ is a diagonal matrix that may differ for each z (Gresele et al., 2021).

Gresele et al. (2021) and Zheng et al. (2022) used this characteristic to find disentangled data generators, but we can again transfer characteristics via faithful encoders: This time we find that searching for an MJ_f with *orthogonal* (instead of disjoint) rows permits post-hoc discovery of concepts. We refer to is property of MJ_f as the *IMA criterion*.

However, as the class of admissible processes has been increased, it is not strong enough to ensure identifiability in the most general case. This is prevented under an additional technical condition on the component magnitudes, which we refer to as *non-equal magnitude ratios* (NEMR). Intuitively, it requires that the magnitudes of the component gradients change non-uniformly between at least two points.

Theorem 3.4 (Identifiability under IMA) Let g adhere to IMA. Let f be a faithful encoder to g. Suppose we have obtained an f' = Mf with a full-rank $M \in \mathbb{R}^{K \times K}$ and orthogonal rows in its Jacobian $MJ_f(g(z)) \coloneqq J_{f'}(g(z))$, i.e, $J_{f'}(g(z))J_{f'}(g(z))^{\top} = \Sigma(z)$ where $\Sigma(z)$ is diagonal and full-rank at two points $z \in \{z_a, z_b\}$. If additionally $\Sigma(z_a)\Sigma(z_b)^{-1}$ has unequal entries its diagonal (NEMR condition), then Me = PSz, where $P \in \mathbb{R}^{K \times K}$ is a permutation and $S \in \mathbb{R}^{K \times K}$ is a scaling matrix.

The constructive proof in App. B.6 can also be condensed into an analytical solution. Alternatively, one can again construct a suitable optimization objective for the IMA criterion, i.e., orthogonal Jacobians. This is achieved by removing the absolute value operation from the arn-operator in Eqn. (1), so that it solely performs a row-wise normalization. In summary, we have established the novel DMA and IMA criteria that allow concept discovery under dependent concepts.

4 EXPERIMENTS

In the following, we perform a battery of experiments of increasing complexity to compare the practical capabilities of approaches for identifiable concept discovery. We start by verifying the theoretical identifiability conditions (Sec. 4.1), then perform evaluation under increasing multi-component correlations for embedding spaces of generative and discriminative models (Sec. 4.2 to 4.4), and finally use a large-scale, discriminatively-trained ResNet50 encoder (Sec. 4.5).

We borrow the DCI metric (Eastwood and Williams, 2018) from disentanglement learning with scores in [0, 1] to measure whether each discovered component predicts precisely one ground-truth component and vice versa. Following Locatello et al. (2020), we report additional metrics with similar results in App. D, along with results on additional datasets and ablations. For reproducibility, each experiment is repeated on five seeds and code is made available upon acceptance. In total, we train and analyze over 300 embedding spaces, requiring about 124 Nvidia RTX2080Ti GPU days. More implementation details are in App. C.2.

4.1 CONFIRMING IDENTIFIABILITY

We first confirm our identifiability guarantees with the analytical solutions. To this end, we implement two realistic synthetic datasets with differentiable generators. This allows computing the closed-form of J_g and deliberately fulfilling or violating the DMA, IMA, and NEMR conditions.

FourBars consists of gray-scale images of four components: Three bars change their colors (black to white) and one bar moves vertically, showing that the image regions affected by each component may change in each image. The plot of J_g in Fig. 3a shows that each component maps to a disjoint image region. This fulfills DMA and thus also IMA. However, all factors have the same gradient magnitudes, making it impossible to find two points with NEMR. According to our theory, we expect DMA optimization to work and IMA to fail. The second dataset, ColorBar, contains a single bar that undergoes realistic changes in color,



(a) FourBars: DMA datasets can be solved by the DMA criterion.



(b) ColorBar: IMA datasets can be solved by the IMA criterion.

Figure 3: Experiments on two synthetic datasets: We confirm our analytic results and show that DMA (a) and IMA (b) cover realistic visual concepts such as colors and translations

width, and its vertical position, see Fig. 3b. It conforms to IMA and NEMR but not DMA. Our proofs indicate that IMA should work and DMA fail. Completing the problem formalization in Sec. 3.1, we compute analytical faithful encoders f for these datasets distorted by a random matrix D. The solutions behave as expected: On FourBars, only the DMA criterion delivers perfectly recovered components (DCI=1) whereas only IMA succeeds on ColorBars.

4.2 CORRELATED COMPONENTS

We now move to the common Shapes3D (Burgess and Kim, 2018) dataset. It shows geometric bodies that vary in their colors, shape, orientation, size, and background totaling six components. Compared to the previous section we train real encoders. We start our analysis where disentanglement learning is no longer possible: When components are correlated. Following Träuble et al. (2021), the dataset is resampled such that two components $z_i, z_j \in [0, 1]$ follow $z_i - z_j \sim \mathcal{N}(0, s^2)$. Lower s results in a stronger correlation where only few pairs of component values co-occur frequently. We choose a moderate correlation of s = 0.4here and three pairs z_i, z_j that are nominal/nominal, nominal/ordinal, and ordinal/ordinal variables. We train four state-of-the-art disentanglement learning VAEs (BetaVAE (Higgins et al., 2017), FactorVAE (Kim and Mnih, 2018), BetaTCVAE (Chen et al., 2018), DipVAE (Kumar et al., 2018)) from a recent study (Locatello et al., 2019) and apply ICA, PCA, and our DMA and IMA discovery methods on their embedding spaces to post-hoc recover the original components. For DMA and IMA, we use the optimizationbased algorithms (Eqn. 1) since they appear more robust to the noisy gradient estimates as demonstrated in App. C.1.

Sec. 4.2 shows the resulting DCI scores. In line with Träuble

Correlated	floor & background		orientation & background		orientation & size	
BetaVAE	0.497 ± 0.03		0.581 ± 0.04		0.491 ± 0.05	
+PCA	0.263 ± 0.03	-47%	0.310 ± 0.02	-47%	0.324 ± 0.04	-34%
+ICA	0.574 ± 0.04	+16%	0.540 ± 0.08	-7%	0.577 ± 0.04	+17%
+Ours (IMA)	0.617 ± 0.02	+24%	0.602 ± 0.05	+3%	0.579 ± 0.03	+18%
+Ours (DMA)	$\textbf{0.641} \pm \textbf{0.03}$	+29%	0.624 ± 0.06	+7%	0.627 ± 0.03	+28%
FactorVAE	0.507 ± 0.11		0.502 ± 0.08		0.712 ± 0.01	
+PCA	0.358 ± 0.07	-29%	0.474 ± 0.05	-5%	0.556 ± 0.03	-22%
+ICA	0.294 ± 0.07	-42%	0.263 ± 0.05	-48%	0.340 ± 0.03	-52%
+Ours (IMA)	0.551 ± 0.04	+9%	0.498 ± 0.03	-1%	0.595 ± 0.05	-16%
+Ours (DMA)	$\textbf{0.584} \pm \textbf{0.05}$	+15%	0.510 ± 0.05	+2%	0.556 ± 0.04	-22%
BetaTCVAE	0.619 ± 0.01		0.613 ± 0.04		0.659 ± 0.01	
+PCA	0.400 ± 0.03	-35%	0.421 ± 0.07	-31%	0.450 ± 0.07	-32%
+ICA	0.540 ± 0.02	-13%	0.497 ± 0.04	-19%	0.627 ± 0.02	-5%
+Ours (IMA)	0.623 ± 0.02	+1%	0.652 ± 0.03	+6%	0.638 ± 0.04	-3%
+Ours (DMA)	$\textbf{0.666} \pm \textbf{0.01}$	+8%	0.664 ± 0.02	+8%	$\textbf{0.748} \pm \textbf{0.03}$	+14%
DipVAE	0.631 ± 0.02		0.652 ± 0.02		0.548 ± 0.04	
+PCA	0.158 ± 0.01	-75%	0.160 ± 0.02	-75%	0.170 ± 0.02	-69%
+ICA	0.630 ± 0.02	-0%	0.651 ± 0.02	-0%	0.542 ± 0.03	-1%
+Ours (IMA)	0.644 ± 0.02	+2%	0.624 ± 0.01	-4%	0.558 ± 0.05	+2%
+Ours (DMA)	$\textbf{0.684} \pm \textbf{0.01}$	+8%	$\textbf{0.679} \pm \textbf{0.01}$	+4%	$\textbf{0.601} \pm \textbf{0.05}$	+10%

Table 2: DMA recovers the components best in 11 out of 12 cases across different models and correlated components of Shapes3D. Mean \pm std. err. of DCI across all components.



Figure 4: DMA discovers directions m that control individual concepts (wall & floor color) of Shapes3D although they are confused in the original embedding space (e_1, e_2) .

et al. (2021), we find that the disentanglement learning VAEs fail to reover the correlated components on their own due to their violated stochastic independence assumption (Fig. 4a). In eleven of the twelve model/correlation pairs, DMA or IMA identify better concepts than the VAE unit axes and the than PCA/ICA components with improvments of up to 29%. This experiment shows that their concept discovery works regardless of (1) the model type and (2) the type of components correlated. On average, DMA delivers better results than IMA (+0.047), despite the generative process of Shapes3D only being roughly IMA or DMA-compliant. This indicates that the DMA criterion might be more robustly optimizable in practice. Fig. 4b visualizes the performance achieved via DMA when traversing the embedding space. It also shows that small DCI differences can mean a significant improvement. This is because (1) the metric is computed across all six components and the strong baselines already identify many concepts and (2) a perfect score of 1.0 is usually not possible due to non-linearly encoded components. We investigate other correlation strengths with similar findings in App. D.3.



Figure 5: DMA and IMA recover the components even under strong and multiple correlations between them. ICA and PCA fail to return better components than the unit axes.

Method	s = 0.1	s = 0.15	s = 0.2	$s = \infty$
unit dirs.	0.238 ± 0.01	0.244 ± 0.01	0.247 ± 0.01	0.286 ± 0.02
PCA	0.238 ± 0.01	0.376 ± 0.03	0.373 ± 0.03	0.343 ± 0.03
ICA	0.409 ± 0.02	0.309 ± 0.02	0.311 ± 0.01	0.652 ± 0.00
(Ours) IMA	0.295 ± 0.01	0.302 ± 0.01	0.333 ± 0.04	0.266 ± 0.12
(Ours) DMA	0.435 ± 0.01	0.411 ± 0.03	0.392 ± 0.02	0.369 ± 0.05

Table 3: Without correlations ($s = \infty$), ICA is able to recover the components of a classification model. Under correlations, DMA works best. Mean \pm std. err. of DCI.

4.3 GAUSSIANITY AND MULTIPLE CORRELATIONS

In this section, we increase the distributional challenges to analyze whether our approaches are as distribution-agnostic as intended. We sample the components of Shapes3D from a (rotationally-symmetric) Gaussian. Additionally, we introduce correlations between multiple components to its covariance matrix. Details on how covariance matrices are constructed are given in App. C.3.

First, we study a single pair of correlated components (floor and background color) with increasing correlation strength ρ . Fig. 5a shows that the BetaVAE handles low correlations well but starts deteriorating from a strength of $\rho > 0.5$, along with ICA. The DCI of our methods is an average constant of +0.145 above the BetaVAE's for $\rho \leq 0.85$. After this, it returns to the underlying BetaVAE's DCI, possibly because the two components collapsed in the BetaVAE's embedding space. For Fig. 5b, we gradually add more moderately correlated ($\rho \approx 0.7$) pairs to the Gaussian's covariance matrix until eventually all components are correlated. Again, our models show a constant benefit over the underlying BetaVAE's DCI curve. This experiment highlights that both DMA and IMA perform well with (1) strong and (2) muliple correlations and (3) Gaussian components.

4.4 DISCRIMINATIVE EMBEDDING SPACES

We highlight that our approach is also applicable to classification models that were trained in a purely discriminative manner, e.g., the feature space of a CNN model. To investigate this setting, we set up an 8-class classification problem



Figure 6: Each component discovered by DMA on CUB correlates with an interpretable ground truth attribute. Images are ordered by their concept scores $(Me)_i$, and the numbers show their ground truth annotated attribute score.

on the Shapes3D dataset, where the combination of the four binarized components object color, wall color (blue/red vs. yellow/green), shape (cylinder vs. cube) and orientation (left vs. right) determines the class as visualized in App. C.4. To make the setting even more realistic, we artificially add labeling noise close to the decision boundary, correlations as in Sec. 4.2, and a small L2-regularizer on the embeddings, keeping them in a reasonable range. We train a discriminative CNN with a K=6-dimensional embedding space.

The discriminative loss leads to a clustered distribution in the embedding space. ICA expectedly works very well in this highly non-Gaussian distribution, when no significant correlations are present which is in line with the result in Theorem 3.4. However, tables turn as we increasingly correlate the floor and background color: Starting at s = 0.2, DMA outperforms ICA and the other methods as can be seen in table Sec. 4.2. While IMA leads to better concepts over the unit directions, it does not reach the level of DMA. Overall this demonstrates that our methods are applicable to purely discriminative embedding spaces and are more robust to high levels of correlations than ICA.

4.5 REAL-WORLD CONCEPT DISCOVERY

Last, we go beyond the traditional benchmarks and perform realistic concept discovery: We analyze the embedding space of a ResNet50 classifier (He et al., 2016) trained on the CUB-200-2011 (Wah et al., 2011) dataset consisting of highresolution images of birds. This amplifies the challenges of the previous sections, i.e., a discriminative space, non-linear component dependencies of varying strengths across multiple components, and a large 512-dimensional embedding space. One restriction of this experiment is that CUB has no data-generating components to compare against, so we cannot report DCI scores. However, we qualitatively show that DMA can deliver interpretable concepts by matching them to annotated attributes of CUB.

We apply DMA and IMA to discover K=30 concepts of which the first two are shown exemplarily in Fig. 6. The images with the highest positive scores on the first component (on the right) consistently show white birds. The other end of the component comprises birds whose primary color is black. This gives a high Spearman rank correlation with the CUB attribute "primary color: white". The second concept is similarly interpretable. To quantify this across all K components, we provide an initial quantitative evaluation based on the Spearman rank correlation between components and attributes in App. D.6. It indicates that ICA and PCA have problems providing such components and the components identified by DMA usually correspond more closely to the attributes. While the construction of further quantitative evaluation schemes goes beyond the scope of this work, these promising results highlight that DMA also works for high-dimensional, real-world datasets.

5 DISCUSSION AND CONCLUSION

Summary. We proposed identifiability as a minimal requirement for concept discovery algorithms. Furthermore, we suggested the two functional paradigms of disjoint and independent mechanisms and proved that they can recover known components in visual embedding spaces. Extensive experiments confirmed that they offer substantial improvements on various generative and discriminative models and remain unaffected by distributional challenges.

Disjoint vs. independent mechanisms. Disjoint mechanisms are grounded in visual compositionality properties, whereas independent mechanisms are grounded in causal principles. In theory, independent mechanisms (IMA) are indisputably more general. Nevertheless, we empirically observed superior results when optimizing disjoint mechanisms (DMA) for complex datasets. We hypothesize that the disjoint loss is more robust to noisy gradient signals and does not suffer from fragile identifiability conditions (e.g., NEMR). Our theoretic results relying on IMA however are more broadly applicable to non-vision tasks and are a valuable starting point for further investigation.

Outlook. We believe our work to be a valuable step towards a rigorous formalization of concept discovery. However, the considered setup can be generalized in the future, for instance to components that are not linearly encoded. This would permit even stronger guarantees. While we have taken a technical perspective, future work is required to investigate the effect of improved concepts on upstream explanations.

References

- Arjun Akula, Shuai Wang, and Song-Chun Zhu. Cocox: Generating conceptual and counterfactual explanations via fault-lines. In *Proceedings of the AAAI Conference* on Artificial Intelligence, volume 34, pages 2594–2601, 2020.
- Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes. *arXiv preprint arXiv:1610.01644*, 2016.
- David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6541–6549, 2017.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- Arianna Bisazza and Clara Tump. The lazy encoder: A finegrained analysis of the role of morphology in neural machine translation. In *Proceedings of the 2018 Conference* on Empirical Methods in Natural Language Processing, pages 2871–2876. Association for Computational Linguistics, 2018.
- Chris Burgess and Hyunjik Kim. 3d shapes dataset. https://github.com/deepmind/3dshapes-dataset/, 2018.
- Christopher P Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. Understanding disentangling in β -VAE. *arXiv preprint arXiv:1804.03599*, 2018.
- Ricky TQ Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disentanglement in variational autoencoders. *Advances in neural information processing systems*, 31, 2018.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020a.
- Zhi Chen, Yijie Bei, and Cynthia Rudin. Concept whitening for interpretable image recognition. *Nature Machine Intelligence*, 2(12):772–782, 2020b.
- Pierre Comon. Independent component analysis, a new concept? *Signal processing*, 36(3):287–314, 1994.
- Jonathan Crabbé and Mihaela van der Schaar. Concept activation regions: A generalized framework for conceptbased explanations. In *Advances in Neural Information Processing Systems*, 2022.

- Cian Eastwood and Christopher KI Williams. A framework for the quantitative evaluation of disentangled representations. In *International Conference on Learning Representations*, 2018.
- Amirata Ghorbani, James Wexler, James Y Zou, and Been Kim. Towards automatic concept-based explanations. In Advances in Neural Information Processing Systems, volume 32, pages 9277–9286, 2019.
- Luigi Gresele, Julius von Kügelgen, Vincent Stimper, Bernhard Schölkopf, and Michel Besserve. Independent mechanism analysis, a new concept? In *Advances in Neural Information Processing Systems*, 2021.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *International Conference on Learning Representations*, 2017.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems, 33:6840–6851, 2020.
- Aapo Hyvärinen and Erkki Oja. A fast fixed-point algorithm for independent component analysis. *Neural computation*, 9(7):1483–1492, 1997.
- Aapo Hyvärinen, Juha Karhunen, and Erkki Oja. Independent component Analysis. John Wiley & Sons, Inc, 2001.
- Aapo Hyvärinen and Petteri Pajunen. Nonlinear independent component analysis: Existence and uniqueness results. *Neural Networks*, 12(3):429–439, 1999. ISSN 0893-6080.
- Aapo Hyvärinen, Hiroaki Sasaki, and Richard Turner. Nonlinear ica using auxiliary variables and generalized contrastive learning. In *The 22nd International Conference* on Artificial Intelligence and Statistics, pages 859–868. PMLR, 2019.
- Ian T Jolliffe. *Principal component analysis*. Springer, 2nd edition, 2002.
- Dmitry Kazhdan, Botty Dimanov, Mateja Jamnik, Pietro Liò, and Adrian Weller. Now you see me (cme): conceptbased model extraction. *AIMLAI workshop at the* 29th ACM International Conference on Information and Knowledge Management (CIKM), 2020.
- Dmitry Kazhdan, Botty Dimanov, Helena Andres Terre, Mateja Jamnik, Pietro Liò, and Adrian Weller. Is disentanglement all you need? comparing concept-based

& disentanglement approaches. *RAI, WeaSul, and Ro*bustML workshops at The Ninth International Conference on Learning Representations 2021, 2021.

- Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvärinen. Variational autoencoders and nonlinear ICA: A unifying framework. In *International Conference* on Artificial Intelligence and Statistics, pages 2207–2217. PMLR, 2020.
- Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *International Conference on Machine Learning*, pages 2668–2677. PMLR, 2018.
- Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *International Conference on Machine Learning*, pages 2649–2658. PMLR, 2018.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. Concept bottleneck models. In *International Conference* on Machine Learning, pages 5338–5348. PMLR, 2020.
- Abhishek Kumar, Prasanna Sattigeri, and Avinash Balakrishnan. Variational inference of disentangled latent concepts from unlabeled observations. In *International Conference on Learning Representations*, 2018.
- Tobias Leemann, Yao Rong, Stefan Kraft, Enkelejda Kasneci, and Gjergji Kasneci. Coherence evaluation of visual concepts with objects and language. In *ICLR2022 Workshop on the Elements of Reasoning: Objects, Structure and Causality*, 2022.
- Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. In *International Conference on Machine Learning*, pages 4114–4124. PMLR, 2019.
- Francesco Locatello, Ben Poole, Gunnar Rätsch, Bernhard Schölkopf, Olivier Bachem, and Michael Tschannen. Weakly-supervised disentanglement without compromises. In *International Conference on Machine Learning*, pages 6348–6359. PMLR, 2020.
- Lukas Muttenthaler, Charles Yang Zheng, Patrick McClure, Robert A. Vandermeulen, Martin N Hebart, and Francisco Pereira. VICE: Variational interpretable concept embeddings. In *Advances in Neural Information Processing Systems*, 2022.

- Bjorn Ommer and Joachim M Buhmann. Learning the compositional nature of visual objects. In 2007 IEEE Conference on Computer Vision and Pattern Recognition, 2007.
- Matthew Painter, Adam Prugel-Bennett, and Jonathon Hare. Linear disentangled representations and unsupervised action estimation. In *Advances in Neural Information Processing Systems*, volume 33, pages 13297–13307. Curran Associates, Inc., 2020.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- Xuanchi Ren, Tao Yang, Yuwang Wang, and Wenjun Zeng. Learning disentangled representation by exploiting pretrained generative models: A contrastive learning view. In *International Conference on Learning Representations*, 2022.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- Frederik Träuble, Elliot Creager, Niki Kilbertus, Francesco Locatello, Andrea Dittadi, Anirudh Goyal, Bernhard Schölkopf, and Stefan Bauer. On disentangled representations learned from correlated data. In *International Conference on Machine Learning*, pages 10401–10412. PMLR, 2021.
- C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, 2011.
- Tao Yang, Xuanchi Ren, Yuwang Wang, Wenjun Zeng, and Nanning Zheng. Towards building a group-based unsupervised representation disentanglement framework. In *International Conference on Learning Representations*, 2021.
- Chih-Kuan Yeh, Been Kim, Sercan O Arik, Chun-Liang Li, Tomas Pfister, and Pradeep Ravikumar. On completenessaware concept-based explanations in deep neural networks. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Yujia Zheng, Ignavier Ng, and Kun Zhang. On the identifiability of nonlinear ICA with unconditional priors. In *ICLR2022 Workshop on the Elements of Reasoning: Objects, Structure and Causality*, 2022.