# DISSECT: Disentangled Simultaneous Explanations via Concept Traversals

**Asma Ghandeharioun**[*]
MIT
asma_gh@mit.edu

**Been Kim**
Google Research
beenkim@google.com

**Chun-Liang Li**
Google Research
chunliang@google.com

**Brendan Jou**
Google Research
bjou@google.com

**Brian Eoff**
Google Research
beoff@google.com

**Rosalind W. Picard**
MIT
picard@media.mit.edu

## Abstract

Explaining deep learning model inferences is a promising venue for scientific understanding, improving safety, uncovering hidden biases, evaluating fairness, and beyond, as argued by many scholars. One of the principal benefits of counterfactual explanations is allowing users to explore "what-if" scenarios through what does not and cannot exist in the data, a quality that many other forms of explanation such as heatmaps and influence functions are inherently incapable of doing. However, most previous work on generative explainability cannot disentangle important concepts effectively, produces unrealistic examples, or fails to retain relevant information. We propose a novel approach, DISSECT, that jointly trains a generator, a discriminator, and a *concept disentangler* to overcome such challenges using little supervision. DISSECT generates Concept Traversals (CTs), defined as a sequence of generated examples with increasing degrees of concepts that influence a classifier's decision. By training a generative model from a classifier's signal, DISSECT offers a way to discover a classifier's inherent "notion" of distinct concepts automatically rather than rely on user-predefined concepts. We show that DISSECT produces CTs that (1) disentangle several concepts, (2) are influential to a classifier's decision and are coupled to its reasoning due to joint training (3), are realistic, (4) preserve relevant information, and (5) are stable across similar inputs. We validate DISSECT on several challenging synthetic and realistic datasets where previous methods fall short of satisfying desirable criteria for interpretability and show that it performs consistently well and better than existing methods. Finally, we present experiments showing applications of DISSECT for detecting potential biases of a classifier and identifying spurious artifacts that impact predictions.

## 1   Introduction

Explanation of the internal inferences of deep learning models remains a challenging problem that many scholars deem promising for improving safety, evaluating fairness, and beyond [10–13]. Many efforts in explainability methods have been working towards providing solutions for this challenging problem. One way to categorize them is by the type of explanations, some post hoc techniques focusing on the importance of individual features, such as saliency maps [1–4], some on importance of individual examples [14–17], some on importance of high-level concepts [18]. There has been active research into the shortcomings of explainability methods (e.g. [19–23]). Scholars have also proposed tests to determine when attention can be used as an explanation [22].

---

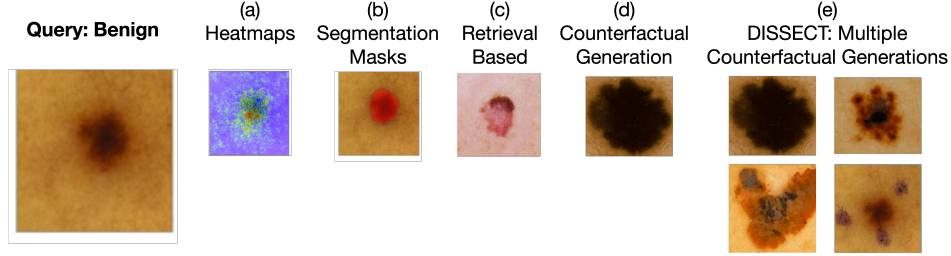[*]Work done in part while at Google Research.

Figure 1: Examples of explainability methods applied to a melanoma classifier. Explanation by (a) heatmaps (e.g. [1–4]), (b) segmentation masks (e.g. [5, 6]), (c) sample retrieval (e.g. [7]), (d) counterfactual generation (e.g. [8, 9]), (e) and multiple counterfactual generations such as DISSECT. Multiple counterfactuals could highlight several different ways that changes in a skin lesion could reveal its malignancy and overcome some of the blind spots of the other approaches. For example, they can demonstrate that large lesions, jagged borders, and asymmetrical shapes lead to melanoma classification. They can also show potential biases of the classifier by revealing that surgical markings can spuriously lead to melanoma classification.

While these methods focus on information *that already exists* in the data, either by weighting features or concepts in training examples or by selecting important training examples, recent progress in generative models [24–28] has lead to another family of explainability methods that provide explanations by *generating new* examples or features [29, 8, 30, 9]. New examples or features can be used to generate *counterfactuals* [31] allowing users to ask: what if this sample were to be classified as the opposite class, and how would it differ? We refer to this line of questioning about the predictor-under-test as "what-if" scenarios. This way of explaining mirrors the way humans reason, justify decisions [32], and learn [33–35]. Additionally, studies have shown that examples are the most preferred means of explanations by users across visual, auditory, and sensor data domains [36].

To illustrate the benefits of counterfactual explanations, consider a dermatology task where an explanation method is used to highlight why a certain sample is classified as benign/malignant (Fig. 1). Explanations like heatmaps, saliency maps, or segmentation masks only provide partial information. Such methods might hint at what is influential within the sample, potentially focusing on the lesion area. However, they cannot show what kind of changes in color, texture, or inflammation could transform the input at hand from benign to malignant. Retrieval-based approaches that provide examples that show a concept are not enough for answering "what-if" questions either. A retrieval-based technique might show input samples of malignant skin lesions that have similarities to a benign lesion in patient A, but from a different patient B, potentially from another body part or even a different skin tone. Such examples do not show what this benign lesion in patient A would have to look like if it were classified as malignant instead. On the other hand, counterfactuals depict *how* to modify the input sample to change its class membership. A counterfactual explanation visualizes what a malignant tumor could look like in terms of potential color or texture, or inflammation changes on the skin. Better yet, multiple counterfactuals could highlight several different ways that changes in a skin lesion could reveal its malignancy. For a classifier that relies on surgical markings [37] as well as meaningful color/texture/size changes to make its decision, a single explanation might fail to reveal this flaw resulting from dependence on spurious features, but multiple distinct explanations shed light on this phenomenon. Multiple explanations are strongly preferred in several other domains, such as educational settings, knowledge discovery, and algorithmic recourse generation [38, 39].

In this work, we propose a generation-based explainability method called DISSECT that generates *Concept Traversals* (CTs)–sequences of generated examples with increasing degrees of concepts' influence on a classifier's decision. While current counterfactual generation techniques fail to satisfy the most consistently agreed-upon properties desired for an explainability method simultaneously [40, 9, 41, 8, 42, 40], DISSECT aims to overcome this challenge. CTs are generated by jointly training a generator, a discriminator, and a CT disentangler to produce examples that (1) express one distinct factor [40] at a time; (2) are influential [9, 41] to a classifier's decision and are coupled to the classifier's reasoning, due to joint training; (3) are realistic [9]; (4) preserve relevant information [8, 42, 40]; and (5) are stable across similar inputs [43, 40, 42]. We compare DISSECT with several baselines, some of which have been optimized for disentanglement, some used extensively for explanation, and some that fall in between. DISSECT is the only technique that performs well across all these dimensions. Other baselines either have a hard time with influence, lack fidelity, generate poor quality and unrealistic samples, or do not disentangle properly. We evaluate DISSECT using

`3D Shapes` [44], `CelebA` [45], and a new synthetic dataset inspired by the challenges faced in the dermatology domain. We show that DISSECT outperforms prior work in addressing all of these challenges. We also discuss this work's applications to detect a classifier's potential biases and identify spurious artifacts using simulated experiments.

This paper makes five main contributions: 1) presents a novel counterfactual explanation approach that manifests several desirable properties outperforming baselines; 2) demonstrates applications through experiments showcasing the effectiveness of this approach for detecting potential biases of a classifier; 3) presents a set of explainability baselines inspired by approaches used for generative disentanglement; 4) translates desired properties commonly referred to in the literature across forms of explanation into measurable quantities for bench-marking and evaluating counterfactual generation approaches; and 5) releases a new synthetic dataset inspired by the challenges faced in the dermatology domain. The code for all the models and metrics is publicly available at `https://github.com/asmadotgh/dissect`.

## 2 Related work

Our method relates to the active research area of post hoc explainability methods. One way to categorize them is by explanation form. While met with criticisms [23, 46], many feature-based explainability methods are shown to be useful [3, 4], which assign a weight to each input feature to indicate their importance in classification. Example-based methods are another popular category [14–17] that instead assign importance weights to individual examples. More recently, concept-based methods have emerged that attribute weights to concepts, i.e., higher-level representations of features [18, 5, 47] such as "long hair". Some of these methods provide multiple explanations [48, 49].

Our work leverages recent progress in generative modeling [50, 51], where the explanation is presented through several conditional generations [52, 53]. Efforts for the "discovery" of concepts are also related to learning disentangled representations [24, 18, 25, 26], which has been shown to be challenging without additional supervision [28]. Recent findings suggest that weak supervision in the following forms could make the problem identifiable: how many factors of variation have changed [54], using labels for validation [55], or using sparse labels during training [55]. While some techniques like Conditional Subspace VAE (CSVAE) [56] began to look into conditional disentanglement by incorporating labels, their performance has not yet reached their unconditional counterparts. Our work uses a new form of weak supervision to improve upon existing methods: the posterior probability and gradient of the classifier-under-test.

Many explainability methods have emerged that use generative models to modify existing examples or generate new examples [57, 58, 30, 8, 29, 9]. Most of these efforts use pre-trained generators and generate a new example by moving along the generator's embedding. Some aim to generate samples that would flip the classifier's decision [29, 58], while others aim to modify particular attribution (e.g., gender) of the image and observe the classifier's decision change [30]. More recent work allows the classifier's predicted probabilities or gradients to flow through the generator during training [9, 8]. However, most assume that there is only one path that crosses the decision boundary, and they generate examples along that path. Our work leverages both counterfactual generation techniques and ideas from the disentanglement literature to generate diverse sets of explanations via multiple distinct paths/concepts influential to the classifier's decision, called Concept Traversals (CT). Each CT is made of a sequence of generated examples that express increasing degrees of a concept.

## 3 Methods

Our key innovation is developing an approach to successfully incorporate the classifier's signal while disentangling concepts to provide multiple potential counterfactuals. We leverage, improve, and build upon several existing methods in parts of our loss function to achieve this goal instead of reinventing the wheel. This section summarizes our final design choices. See §A.1 and §A.2 for more details.

**Notation.** Without loss of generality, we consider a binary classifier $f\colon X \to Y = \{-1, 1\}$ such that $f(x) = p(y = 1|x)$ where $x \sim \mathcal{P}_X$. We want to find $K$ latent concepts that contribute to the decision-making of $f$. Given $x$, $\alpha \in [0, 1]$, and $k$, we want to generate an image, $\bar{x}$, by perturbing the latent concept $k$, such that the posterior probability $f(\bar{x}) = \alpha$. In addition, when

conditioning on $x$ and a concept $k$, by changing $\alpha$, we hope for a smooth change in the output image $\bar{x}$. This smoothness resembles slightly changing the degree of a concept and aids in interpretation. Putting it together, we desire a generic generation function that generates $\bar{x}$, defined as $\mathcal{G}(x, \alpha, k; f) : X \times [0, 1] \times \{0, 1, \dots, K-1\} \to X$, where $f(\mathcal{G}(x, \alpha, k; f)) \approx \alpha$. The generation function $\mathcal{G}$ conditions on $x$ and manipulates the concept $k$ such that a monotonic change in the $f(\mathcal{G}(x, 0, k)), \dots, f(\mathcal{G}(x, 1, k))$ is achieved.

### 3.1 Encoder-decoder architecture

We realize the generic (conditional) generation process using an encoder-decoder framework. The input $x$, is encoded into an embedding $E(x)$, before feeding into the generator $G$, which serves as a decoder. In addition to $E(x)$, we have a conditional input $c(\alpha, k)$ for the $k^{\text{th}}$ latent concept we are going to manipulate with level $\alpha$. The generative process $\mathcal{G}$ can be implemented by $G(E(x), c(\alpha, k))$, where we can manipulate $c$ for interpretation via conditional generation.

There are many advances in generative models, which can be adopted to train $G(E(x), c(\alpha, k))$. For example, Variational Autoencoders (VAE) explicitly designed for disentanglement, such as $\beta$-VAE [24], Annealed-VAE [59], $\beta$-TCVAE [60], and DIPVAE [61]. More advanced approaches include Conditional subspace VAE (CSVAE) [56]. However, VAE-based approaches often generate blurry images with some exceptions (e.g., adopting customized architecture designs [62, 63]). In contrast, GANs [51] tend to generate high-quality images. Therefore, we consider an encoder-decoder design with GANs [64–69] to have flexible controllability and high quality generation for explainability.

There are different design choices for realizing the conditioning code $c$. A straightforward option is multiplying a $K$-dimensional one-hot vector, which specifies the desired concept $k$, with the probability $\alpha$. However, conditioning on continuous variables is usually challenging for conditional GANs compared with their discretized counterparts [53]. Therefore, we discretize $\alpha$ into $[0, 1/N, \dots, 1]$, and parametrize $c$ as a binary matrix $[0, 1]^{(N+1) \times K}$ for $N + 1$ discrete values and $K$ concepts. Therefore, $c(\alpha, k)$ denotes that we set the elements $(m, k)$ to be 1, where $m$ is an integer and $m/N \le \alpha < (m+1)/N$. See § A.2 for more details.

We introduce our loss function, which contains different components for our desired properties.

**Interpreting classifiers.** The main difference between the proposed model and the classic generative models is the ability to produce generator outputs aligned with the classifier being interpreted. That is,

$$\mathcal{L}_f(G) = \texttt{LogLoss}\Big(\alpha, f\big(G\left(x, c(\alpha, k)\right)\big)\Big),$$

where $\texttt{LogLoss}$ is the log loss for binary classification.

**Reconstruction.** Under an encoder-decoder framework, we require the model to be able to reconstruct the data. Given $x \sim \mathcal{P}_X$, and a classifier $f$, the ground-truth $\alpha$ for all concepts is $f(x)$. We adopt $\ell_1$ reconstruction loss as used in [66–69]

$$\mathcal{L}_{\text{rec}}(G) = \frac{1}{K} \sum_{k=0}^{K-1} \left\| x - G\big(x, c\left(f(x), k\right)\big) \right\|_1.$$

**Cycle consistency.** In addition to the standard reconstruction, we regularize the model with cycle consistency [69]. Define $\bar{x}_{\alpha, k} = G(x, c(\alpha, k))$, which is the perturbed version of $x$ by manipulating the concept $k$ to be level $\alpha$. The only difference between $x$ and $\bar{x}$ is the difference in levels of concept $k$. Therefore, if we condition on $\bar{x}_{\alpha, k}$, we should be able to reconstruct $x$ by changing concept $k$ from level $\alpha$ to $f(x)$, that is

$$\mathcal{L}_{\text{cyc}}(G) = \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}_\alpha \left[ \left\| x - G\big(\bar{x}_{\alpha, k}, c\left(f(x), k\right)\big) \right\|_1 \right].$$

**GAN loss.** As in previous work [64–68] using GANs to enhance generation quality, we use GAN loss to compare the distribution between $x \sim \mathcal{P}_X$ and the generated data with uniformly sampled

concept $k$ and level $\alpha$. One can use any GAN loss [51, 70–74]. We use spectral normalization [74] with its hinge loss variant [53] following [9], where

$$\mathcal{L}_{\text{cGAN}}(D) = -\mathbb{E}_{x \sim \mathcal{P}_X}\left[ \min\left(0, -1 + D(x)\right) \right] - \mathbb{E}_{x \sim \mathcal{P}_X} \mathbb{E}_{\alpha,k}\left[ \min\left(0, -1 - D\big(G\left(x, c(\alpha, k)\right)\big)\right) \right],$$

and

$$\mathcal{L}_{\text{cGAN}}(G) = -\mathbb{E}_{x \sim \mathcal{P}_X} \mathbb{E}_{\alpha,k}\left[ D\big(G\left(x, c(\alpha, k)\right)\big) \right].$$

### 3.2 Enforcing disentanglement

Ideally, we want the $K$ concepts to capture distinctively different qualities influential to the classifier being explained. Though $G$ is capable of expressing $K$ different concepts, the above formulation does not enforce distinctness among them. To address this challenge, we introduce a *concept disentangler* $R$ inspired by the disentanglement literature [24, 59–61, 56, 26, 75]. Given a pair of images $(x, x')$, $R$ tries to predict which concept $k$, where $k \in \{0, \dots, K-1\}$ has perturbed query $x$ to produce $x'$. To formalize this, let $\bar{x}_{k,\alpha} = G(x, c(\alpha, k))$ and $\tilde{x}_k = G(\bar{x}_{k,\alpha}, c(f(x), k))$. The loss function is

$$\mathcal{L}_r(G, R) = \mathbb{CE}(e_k, R(x, \bar{x}_k)) + \mathbb{CE}(e_k, R(\bar{x}_k, \tilde{x}_k)),$$

where $\mathbb{CE}$ is the cross entropy and $e_k$ is a one-hot vector of size $K$ for the $k^{\text{th}}$ concept. By jointly optimizing $R$ and $G$, we penalize the generator $G$ if any concept shares similar information with others. In summary, the overall objective function of our method is:

$$\min_{G,R} \max_{D} \lambda_{\text{cGAN}} \mathcal{L}_{\text{cGAN}}(D, G) + \lambda_f \mathcal{L}_f(D, G) + \lambda_{\text{rec}} \mathcal{L}_{\text{rec}}(G) + \lambda_{\text{rec}} \mathcal{L}_{\text{cyc}}(G) + \lambda_r \mathcal{L}_r(G, R).$$

We call our approach DISSECT as it disentangles simultaneous explanations via concept traversals.

## 4 Experiments

### 4.1 Datasets

**3D Shapes.** We first use 3D Shapes [44], a synthetic dataset available under Apache License 2.0 composed of 480K 3D shapes procedurally generated from 6 ground-truth factors of variation. These factors are floor hue, wall hue, object hue, scale, shape, and orientation. Note that this dataset is used purely for validation and demonstration purposes due to the controllability of all these factors.

**SynthDerm.** Second, we create a new dataset, SynthDerm (Fig. 2). Real-world characteristics of melanoma skin lesions in dermatology settings inspire the generation of this dataset [76]. These characteristics include whether the lesion is asymmetrical, its border is irregular or jagged, is unevenly colored, has a diameter more than 0.25 inches, or is evolving in size, shape, or color over time. These qualities are usually referred to as the ABCDE of melanoma [77]. We generate SynthDerm algorithmically by varying several factors: skin tone, lesion shape, lesion size, lesion location (vertical and horizontal), and whether there are surgical markings present. We randomly assign one of the following to the lesion shape: round, asymmetrical, with jagged borders, or multi-colored (two different shades of colors overlaid with salt-and-pepper noise). For skin tone values, we simulate Fitzpatrick ratings [78]. Fitzpatrick scale is a commonly used approach to classify the skin by its reaction to sunlight exposure modulated by the density of melanin pigments in the skin. This rating has six values, where 1 represents skin that always burns (lowest melanin) and 6 represents skin that never burns in sunlight (highest melanin). For our synthetic generation, we consider six base skin tones that similarly resemble different amounts of melanin. We also add a small amount of random noise to the base color to add further variety. Overall, SynthDerm includes more than 2,600 images of size 64x64. We have made this dataset publicly available at https://affect.media.mit.edu/dissect/synthderm.

**CelebA.** We also include the CelebA dataset [45] that is available for non-commercial research purposes only. This dataset includes "identities [of celebrities] collected from the Internet" [79], is realistic, and closely resembles real-world settings where the attributes are nuanced and not mutually independent . CelebA includes more than 200K images with 40 annotated face attributes, such as smiling, hair color, and bangs. We could not find any information about the consent procedure.
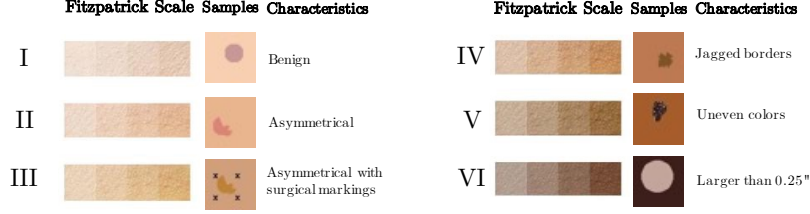
| Fitzpatrick Scale | Samples | Characteristics | | Fitzpatrick Scale | Samples | Characteristics |
|---|---|---|---|---|---|---|
| I | | Benign | | IV | | Jagged borders |
| II | | Asymmetrical | | V | | Uneven colors |
| III | | Asymmetrical with surgical markings | | VI | | Larger than 0.25" |

Figure 2: Illustration of `SynthDerm` dataset that we algorithmically generated. Fitzpatrick scale of skin classification based on melanin density and corresponding examples in the dataset are visualized.

## 4.2 Baselines

We consider several baselines. First, we modify a set of VAEs explicitly designed for disentanglement [24, 59–61] to incorporate the classifier's signal during their training processes and encourage the generative model to learn latent dimensions that influence the classifier, i.e., learning *Influential* CTs. We refer to these baselines with a **-mod** postfix, e.g., $\beta$**-VAE-mod**. Second, we include **CSVAE** [56] that aims to solve unsupervised learning of features associated with a specific label using a low-dimensional latent subspace that can be independently manipulated. Third, we include Explanation by Progressive Exaggeration (**EPE**) [9], a GAN-based approach that learns to generate one series of counterfactual and realistic samples that change the prediction of $f$, given data and the classifier's signal. We introduce **EPE-mod** by modifying EPE to allow learning multiple pathways by making the generator conditional on another variable: the CT dimension (see more details in §A.1).

## 4.3 Evaluation metrics

To evaluate the quality of the discovered CTs, we consider several measures that formalize *Importance* [9, 41], *Realism* [9], *Distinctness* [40], *Substitutability* [42, 40, 8, 80], and *Stability* [43, 40, 42], which commonly appear as desired qualities in the explainability literature.

**Importance.** Explanations should produce the desired outcome from the black-box classifier $f$. Previous work has referred to this quality using different names, such as importance [5], compatibility with classifier [9], and classification model consistency [41]. While most previous methods have relied on visual inspection, we introduce a quantitative metric to measure the gradual increase of the target class's posterior probability through a CT. Notably, we compute the correlation between $\alpha$ and $f(I(x, \alpha, k; f))$ introduced in § 3. For brevity, we refer to $f(I(x, \alpha, k; f))$ as $f(\bar{x})$ in the remainder of the paper. We also report the mean-squared error and the Kullback–Leibler (KL) divergence between $\alpha$ and $f(\bar{x})$. We also calculate the performance of the black-box classifier $f$ on the counterfactual explanations. Specifically, we replace the test set of real images with their generated counterfactual explanations and quantify the performance of the pre-trained black-box classifier $f$ on the counterfactual test set. Better performance suggests that counterfactual samples are compatible with $f$ and lie on the correct side of the classifier's boundary.

**Realism.** We need the generated samples that form a CT to look realistic to enable users/humans to identify concepts they represent. It means the counterfactual explanations should lie on the data manifold. This quality has been referred to as realism or data consistency [9]. Inspired by [81], we train a post hoc classifier that predicts whether a sample is real or generated. Although its objective is identical to that of the discriminator in our architecture, it is essential to do this step post hoc and independent from the training procedure because relying on the discriminator's accuracy in an adversarial training framework can be misleading [81].

**Distinctness.** Another desirable quality for explanations is to represent inputs with non-overlapping concepts, often referred to as diversity [40]. Others have suggested similar properties such as coherency, meaning examples of a concept should be similar but different from examples of other concepts [5]. To quantify this in a counterfactual generation setup, we introduce a distinctness measure capturing the performance of a secondary classifier that distinguishes between CTs. Mainly, we train a classifier post hoc that given a query image $x$ and a generated image $x'$ and $K$ number of CTs, predicts $CT_k$ has perturbed $x$ to produce $x'$, $k \in \{1, 2, \dots, K\}$. This classifier is agnostic to our model and only uses its pair of input images.
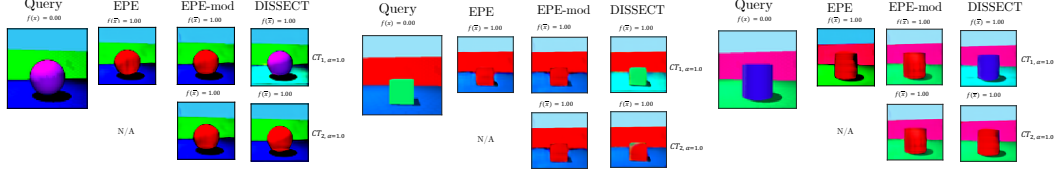
Figure 3: Examples from 3D Shapes. EPE and EPE-mod converge to finding the same concept, despite EPE-mod's ability to express multiple pathways to switch classifier outcomes from False to True. DISSECT discovers the two ground-truth concepts: $CT_1$ flips the floor color to cyan and $CT_2$ flips the shape color to red.

Table 1: Quantitative results on 3D Shapes. DISSECT performs significantly better or on par with the strongest baselines in all evaluation categories. Modified VAE variants perform poorly in terms of *Importance*, worse than CSVAE, and significantly worse than EPE, EPE-mod, and DISSECT. CSVAE and other VAE variants do not produce high-quality images, thus have poor *Realism* scores; meanwhile, EPE, EPE-mod, and DISSECT generate realistic samples indistinguishable from real images. While the aggregated metrics for *Importance* are useful for discarding VAE baselines with poor performance, they do not show a consistent order across EPE, EPE-mod, and DISSECT. Our approach greatly improves *Distinctness*, especially compared to EPE-mod. EPE is inherently incapable of doing this, and the extension EPE-mod does, but poorly. For the *Substitutability* scores, note the classifier's precision, recall, and accuracy when training on actual data is 100%. (*Certain VAE methods only generate samples with $f(\bar{x}) = 0.0$. Correlation with a constant value is undefined.)

| | Importance | | | | | | | Realism | | | Distinctness | | | | | Substitutability | | | Stability | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | ↑R | ↑ρ | ↓KL | ↓MSE | ↑CF Acc | ↑CF Prec | ↑CF Rec | ↓Acc | ↓Prec | ↓Rec | ↑Acc | ↑Prec (micro) | ↑Prec (macro) | ↑Rec (micro) | ↑Rec (macro) | ↑Acc Sub | ↑Prec Sub | ↑Rec Sub | ↓CF MSE | ↓Prob JSD |
| VAE-mod | 0.1 | 0.3 | inf | 0.42 | 50.0 | 0 | 0 | 94.9 | 92.1 | 99.6 | 86.3 | 95.1 | 95.3 | 80.6 | 80.6 | 14.7 | 14.2 | 69.4 | 0.151 | **0.0000** |
| β-VAE-mod | 0.0 | 0.0 | inf | 0.42 | 50.0 | 0 | 0 | 99.5 | 99.0 | 100 | 90.7 | 92.2 | 92.2 | 91.0 | 91.0 | 42.5 | 8.2 | 19.8 | 0.197 | **0.0000** |
| Annealed-VAE-mod | N/A* | N/A* | inf | 0.42 | 50.0 | 0 | 0 | 100 | 100 | 100 | 34.0 | 53.2 | 49.2 | 1.20 | 1.2 | 35.8 | 14.4 | 48.1 | 0.175 | **0.0000** |
| DIPVAE-mod | 0.1 | 0.5 | inf | 0.41 | 52.3 | 100 | 4.6 | 100 | 100 | 100 | 97.2 | 96.7 | 96.9 | 96.5 | 96.5 | 19.0 | 19.0 | 100 | 0.127 | 0.0001 |
| βTCVAE-mod | 0.5 | 0.7 | inf | 0.42 | 50.0 | 0 | 0 | 100 | 100 | 100 | **100** | **100** | **100** | **100** | **100** | 36.4 | 21.3 | 87.3 | 0.143 | **0.0000** |
| CSVAE | 0.3 | 0.3 | 5.5 | 0.28 | 64.6 | 100 | 29.3 | 100 | 100 | 100 | 71.4 | 74.7 | 76.4 | 87.2 | 87.2 | 47.0 | 23.8 | 81.3 | 19.544 | 0.0274 |
| EPE | 0.8 | **0.8** | 1.54 | 0.09 | 98.4 | 100 | 96.7 | 50.1 | **0** | **0** | - | - | - | - | - | 99.2 | 95.9 | **100** | 0.134 | 0.0004 |
| EPE-mod | **0.9** | 0.7 | 2.2 | **0.08** | **99.7** | 100 | **99.4** | **49.3** | **0** | **0** | 45.3 | 49.6 | 49.8 | 30.3 | 30.3 | 91.0 | **100** | 52.5 | 0.128 | 0.0002 |
| **DISSECT** | 0.8 | **0.8** | 1.61 | **0.08** | 98.7 | 100 | 97.5 | **49.3** | **0** | **0** | **100** | **100** | **100** | **100** | **100** | **100** | 99.7 | **100** | **0.102** | 0.0003 |

**Substitutability.** The representation of a sample in terms of concepts should preserve relevant information [42, 40]. Previous work has formalized this quality for counterfactual generation contexts through a proxy called substitutability [8] . Substitutability measures an external classifier's performance on real data when it is trained using only synthetic images.[2] A higher substitutability score suggests that explanations have retained relevant information and are of high quality.

**Stability.** Explanations should be coherent for similar inputs, a quality known as stability [43, 40, 42]. To quantify stability in counterfactual explanations, we augment the test set with additive random noise to each sample $x$ and produce $S$ randomly augmented copies $x_i''$. Then, we generate counterfactual explanations $\bar{x}$ and $\bar{x}_i''$, respectively. We calculate the mean-squared difference between counterfactual images $\bar{x}$ and $\bar{x}_i''$ and the resulting Jensen Shannon distance (JSD) between the predicted probabilities $f(\bar{x})$ and $f(\bar{x}_i'')$.

### 4.4 Case study I: validating the qualities of concept traversals

Considering 3D Shapes [44], we define an image as "colored correctly" if the shape hue is red or the floor hue is cyan. We train a classifier to detect whether a sample image is "colored correctly" or not. In this synthetic experiment, only these two independent factors contribute to the decision of this classifier. Given a not "colored correctly" query, we would like the algorithm to find a CT related to the shape color and another CT associated with the floor color–two different pathways that lead to switching the classifier outcome for that sample.[3]

Tab. 1 summarizes the quantitative results on 3D Shapes. Most VAE variants perform poorly in terms of *Importance*, CSVAE performs slightly better, and EPE, EPE-mod, and DISSECT perform best.

---

[2]This metric has been used in other contexts outside of explainability and has been called Classification Accuracy Score (CAS) [80]. CAS is more broadly applicable than Fréchet Inception Distance [82] and Inception Score [83] that are only useful for evaluating GAN models.

[3]In this scenario, these two ground-truth concepts do not directly apply to switching the classifier outcome from True to False. For example, if an image has a red shape and a cyan floor, both of these colors need to be changed to switch the classification outcome. We still observe that DISSECT finds CTs that change different combinations of colors, but the baseline methods converge to the same CT. See §A.5 for more details.
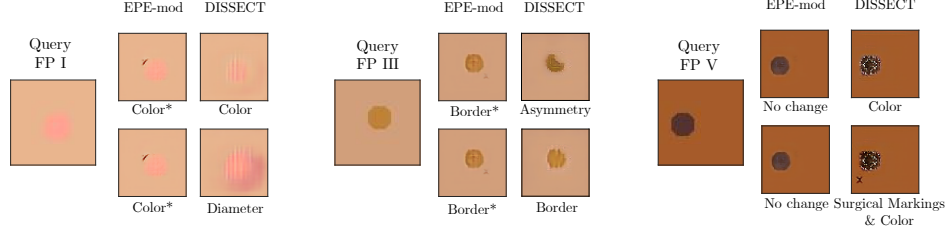
Figure 4: Examples from `SynthDerm` comparing DISSECT with the strongest baseline, EPE-mod. We illustrate three queries with different Fitzpatrick ratings [78] and visualize the two most prominent concepts for each technique. We observe that EPE-mod converges on a single concept that only vaguely represents meaningful ground-truth concepts. However, DISSECT successfully finds concepts describing asymmetrical shapes, jagged borders, and uneven colors that align with the ABCDE of melanoma [77]. DISSECT also identifies concepts for surgical markings that spuriously impact the classifier's decisions.

Table 2: Quantitative results on `SynthDerm`. DISSECT performs consistently best in all categories. For anchoring the *Substitutability* scores, note that the precision, recall, and accuracy of the classifier when training on actual data is 97.7%, 100.0%, and 95.4%, respectively.

| | Importance | | | | | | | Realism | | | Distinctness | | | | | Substitutability | | | Stability | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | ↑R | ↑ρ | ↓KL | ↓MSE | ↑CF Acc | ↑CF Prec | ↑CF Rec | ↓Acc | ↓Prec | ↓Rec | ↑Acc | ↑Prec (micro) | ↑Prec (macro) | ↑Rec (micro) | ↑Rec (macro) | ↑Acc Sub | ↑Prec Sub | ↑Rec Sub | ↓CF MSE | ↓Prob JSD |
| CSVAE | 0.25 | 0.64 | 1.78 | 0.12 | 86.7 | 43.9 | 5.2 | 54.6 | 69.9 | 16.3 | 85.1 | 0 | 0 | 0 | 0 | 29.9 | 36.8 | 55.4 | 2.318 | 0.006 |
| EPE | 0.87 | 0.23 | inf | 0.03 | 80.7 | 55.1 | 86.9 | 50.1 | **0** | **0** | - | - | - | - | - | 74.4 | 83.8 | 60.7 | **0.111** | **0.001** |
| EPE-mod | 0.81 | 0.73 | 0.92 | 0.04 | 95.3 | 83.9 | 79.5 | **50.0** | **0** | **0** | 85.1 | 71.1 | 26.3 | 0.7 | 0.7 | 74.3 | 83.5 | 60.7 | 0.239 | 0.002 |
| **DISSECT** | **0.92** | **0.75** | **0.35** | **0.02** | **97.8** | **92.3** | **91.1** | **50.0** | **0** | **0** | **96.0** | **96.5** | **96.5** | **74.6** | **74.6** | **81.0** | **97.2** | **64.0** | 0.338 | **0.001** |

Our results suggest that DISSECT performs similarly to EPE that has been geared explicitly toward exhibiting *Importance* and its extension, EPE-mod. Additionally, DISSECT still keeps *Realism* intact. Also, it notably improves the *Distinctness* of CTs compared to relevant baselines.

Fig. 3 shows the examples illustrating the qualitative results for EPE, EPE-mod, and DISSECT. Our results reveal that EPE converges to finding only one of these concepts. Similarly, both CTs generated by EPE-mod converge to finding the same concept, despite being given the capability to explore two pathways to switch the classifier outcome. However, DISSECT finds the two distinct ground-truth concepts through its two generated CTs. For brevity, three sample queries are visualized (See more in §A.5).

## 4.5 Case study II: identifying spurious artifacts

A "high performance" model could learn to make its decisions based on irrelevant features that only happen to correlate with the desired outcome, known as label leakage [84]. One of the applications of DISSECT is to uncover such spurious concepts and allow probing a black-box classifier. Motivated by real-world examples that revealed classifier dependency on surgical markings in identifying melanoma [37], we design this experiment. Given the synthetic nature of `SynthDerm` and how it has been designed based on real-world characteristics of melanoma [76, 77], each sample has a deterministic label of melanoma or benign. If the image is asymmetrical, has jagged borders, has different colors represented by salt-and-pepper noise, or has a large diameter (i.e., does not fit in a 40×40 square), the sample is melanoma. Otherwise, the image represents the benign class. Similar to in-situ dermatology images, melanoma samples have surgical markings more frequently than benign samples. We train a classifier to detect whether a sample image is melanoma or a benign lesion.

Given a benign query, we would like to produce counterfactual explanations that depict *how* to modify the input sample to change its class membership. We want DISSECT to find CTs that disentangle meaningful characteristics of melanoma identification in terms of color, texture, or shape [77], and identify potential spurious artifacts that impact the classifier's predictions.

Tab. 2 summarizes the quantitative results on `SynthDerm`. Our method performs consistently well across all the metrics, significantly boosting *Distinctness* and *Substitutability* scores and making meaningful improvements on *Importance* scores. Our approach has higher performance compared to EPE-mod and EPE baselines and substantially improves upon CSVAE. Our method's high *Distinctness* and *Substitutability* scores show that DISSECT covers the landscape of potential concepts very well and retains the variety seen in real images strongly better than all the other baselines.
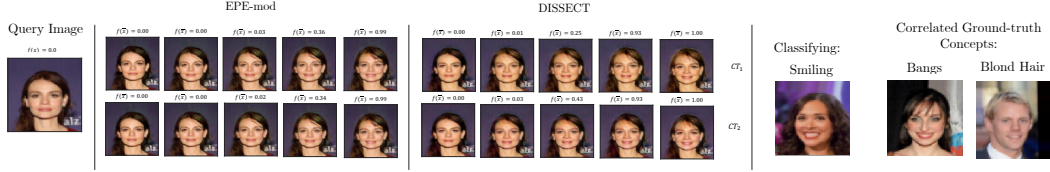
Figure 5: Examples from `CelebA`. A biased classifier has been trained to predict smile probability, where a training dataset was sub-sampled so that smiling co-occurs only with "bangs" and "blond hair" attributes. EPE does not support multiple CTs. EPE-mod converges on the same concept, despite having the ability to express various pathways to change $f(\bar{x})$ through $CT_1$ and $CT_2$. However, DISSECT discovers distinct pathways: $CT_1$ mainly changes hair color to blond, and $CT_2$ does not alter hair color but focuses on hairstyle and tries to add bangs. DISSECT identifies two otherwise hidden biases.

Table 3: Quantitative results on `CelebA`. DISSECT performs better than or on par with the baselines in all categories. Notably, DISSECT greatly improves the *Distinctness* of CTs and achieves a higher *Realism* score, suggesting disentangling CTs does not diminish the quality of generated images and may even improve them. For anchoring *Substitutability* scores, note that the classifier's precision, recall, and accuracy when training on actual data is 95.4%, 98.6%, and 92.7%, respectively.

| | Importance | | | | | | | Realism | | | Distinctness | | | | | Substitutability | | | Stability | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | ↑R | ↑ρ | ↓KL | ↓MSE | ↑CF Acc | ↑CF Prec | ↑CF Rec | ↓Acc | ↓Prec | ↓Rec | ↑Acc | ↑Prec (micro) | ↑Prec (macro) | ↑Rec (micro) | ↑Rec (macro) | ↑Acc Sub | ↑Prec Sub | ↑Rec Sub | ↓CF MSE | ↓Prob JSD |
| CSVAE | 0.00 | 0.00 | 1.25 | 0.284 | 50.2 | 50.2 | 34.1 | 99.7 | 100 | 99.5 | 15.1 | 0.0 | 0.0 | 0.0 | 0.0 | 52.8 | 53.0 | **98.6** | 23.722 | 0.039 |
| EPE | 0.85 | **0.91** | 0.28 | 0.060 | 99.2 | **99.9** | 98.5 | 49.9 | 32.0 | 0.3 | - | - | - | - | - | **93.0** | 95.4 | 91.2 | **0.411** | **0.003** |
| EPE-mod | **0.86** | 0.90 | 0.21 | 0.048 | **99.5** | 99.7 | **99.2** | 49.3 | 33.3 | 0.1 | 18.7 | 50.0 | 50.1 | 15.2 | 15.2 | 91.8 | 94.8 | 89.4 | 0.446 | 0.004 |
| **DISSECT** | 0.84 | 0.88 | **0.19** | **0.047** | 99.2 | 99.8 | 98.5 | **49.2** | **0.0** | **0.0** | **95.0** | **98.0** | **98.1** | **96.1** | **96.1** | 91.9 | **96.9** | 87.6 | 0.567 | 0.005 |

Fig. 4 illustrates a few examples to showcase DISSECT's improvements over the strongest baseline, EPE-mod. EPE-mod converges to finding a single concept that only vaguely represents meaningful ground-truth concepts. However, DISSECT successfully finds concepts describing asymmetrical shapes, jagged borders, and uneven colors that align with ABCDE of melanoma [77]. DISSECT also identifies surgical markings as a spurious concept that impacts the classifier's decisions. Overall, the qualitative results show that DISSECT uncovers several critical blind spots of the baseline techniques.

## 4.6 Case study III: identifying biases

Another potential use case of DISSECT is to identify biases that might need to be rectified. Since our approach does not depend on predefined user concepts, it may help discover unkwown biases. We design an experiment to test DISSECT in such a setting. We sub-sample `CelebA` to create a training dataset such that smiling correlates with "blond hair" and "bangs" attributes. In particular, positive samples either have blond hair or have bangs, and negative examples are all dark-haired and do not have bangs. We use this dataset to train a purposefully biased classifier. We employ DISSECT to generate two CTs. Fig. 5 shows the qualitative results, which depict that DISSECT automatically discovers the two biases, which other techniques fail to do. Tab. 3 summarizes the quantitative results that replicate our finding from Tab. 1 in § 4.4 and Tab. 2 in § 4.5 in a real-world dataset, confirming that DISSECT quantitatively outperforms all the other baselines in *Distinctness* without negatively impacting *Importance* or *Realism*.

## 5 Concluding remarks

We present DISSECT that successfully finds multiple distinct concepts by generating a series of realistic-looking, counterfactually generated samples that gradually traverse a classifier's decision boundary. We validate our approach by experimental results, both quantitatively and qualitatively showing significant improvement over prior methods. Mirroring natural human reasoning for explainability [32–35], the new method provides additional checks-and-balances for practitioners to probe a trained model prior to deployment, especially in high-stakes tasks such as medical decision making. DISSECT helps identify how well the model-under-test reflects practitioners' domain knowledge and whether the model exhibits biases that might need to be rectified. Since this method does not depend on predefined user concepts, it may also help discover biases that were not anticipated.

One avenue for future explainability work is extending earlier theories [28, 54] that obtain disentanglement guarantees. While we provide an extensive list of qualitative examples in the Appendix,

further confirmation from human-subject studies to validate that CTs exhibit semantically meaningful attributes could strengthen our findings. We emphasize that our proposed method does not guarantee finding all the biases of a classifier, nor ensures semantic meaningfulness across all found concepts. Additionally, it should not replace procedures for promoting transparency in model evaluation such as model cards [85]; instead, we recommend it be used in tandem. We also acknowledge the limitations for Fitzpatrick ratings used in `SynthDerm` dataset in capturing variability, especially for darker skin tones [86, 87], and would like to consider other ratings in the future.

## Acknowledgments and Disclosure of Funding

## References

[1] Dumitru Erhan, Yoshua Bengio, Aaron Courville, and Pascal Vincent. Visualizing higher-layer features of a deep network. *University of Montreal*, 1341(3):1, 2009.

[2] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. SmoothGrad: Removing noise by adding noise. *ICML Workshop on Visualization for Deep Learning*, 2017.

[3] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, volume 70, pages 3319–3328, 2017.

[4] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Advances in neural information processing systems*, 2017.

[5] Amirata Ghorbani, James Wexler, James Y Zou, and Been Kim. Towards automatic concept-based explanations. In *Advances in Neural Information Processing Systems*, volume 32, 2019.

[6] Alberto Santamaria-Pang, James Kubricht, Aritra Chowdhury, Chitresh Bhushan, and Peter Tu. Towards emergent language symbolic semantic segmentation and model interpretability. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 326–334. Springer, 2020.

[7] Wilson Silva, Alexander Poellinger, Jaime S Cardoso, and Mauricio Reyes. Interpretability-guided content-based medical image retrieval. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 305–314. Springer, 2020.

[8] Pouya Samangouei, Ardavan Saeedi, Liam Nakagawa, and Nathan Silberman. ExplainGAN: Model explanation via decision boundary crossing transformations. In *European Conference on Computer Vision*, pages 666–681, 2018.

[9] Sumedha Singla, Brian Pollack, Junxiang Chen, and Kayhan Batmanghelich. Explanation by progressive exaggeration. In *International Conference on Learning Representations*, 2020. URL `https://openreview.net/forum?id=H1xFWgrFPS`.

[10] Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*, 2017.

[11] Alyssa Glass, Deborah L McGuinness, and Michael Wolverton. Toward establishing trust in adaptive agents. In *Proceedings of the 13th international conference on Intelligent user interfaces*, pages 227–236, 2008.

[12] Henriette Cramer, Jean Garcia-Gathright, Aaron Springer, and Sravana Reddy. Assessing and addressing algorithmic bias in practice. *Interactions*, 25(6):58–63, 2018.

[13] Charles T Marx, Richard Lanas Phillips, Sorelle A Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian. Disentangling influence: Using disentangled representations to audit model predictions. *arXiv preprint arXiv:1906.08652*, 2019.

[14] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *International Conference on Machine Learning*, volume 70, pages 1885–1894, 2017.

[15] Been Kim, Cynthia Rudin, and Julie A Shah. The bayesian case model: A generative approach for case-based reasoning and prototype classification. In *Advances in neural information processing systems*, pages 1952–1960, 2014.

[16] Chih-Kuan Yeh, Joon Kim, Ian En-Hsu Yen, and Pradeep K Ravikumar. Representer point selection for explaining deep neural networks. In *Advances in neural information processing systems*, 2018.

[17] Rajiv Khanna, Been Kim, Joydeep Ghosh, and Sanmi Koyejo. Interpreting black box predictions using Fisher kernels. In *International Conference on Artificial Intelligence and Statistics*, pages 3382–3390, 2019.

[18] 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 2673–2682, 2018.

[19] Danish Pruthi, Mansi Gupta, Bhuwan Dhingra, Graham Neubig, and Zachary C Lipton. Learning to deceive with attention-based explanations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4782–4793, 2020.

[20] Sofia Serrano and Noah A Smith. Is attention interpretable? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2931–2951, 2019.

[21] Sarthak Jain and Byron C Wallace. Attention is not explanation. In *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, volume 1, pages 3543–3556, 2019.

[22] Sarah Wiegreffe and Yuval Pinter. Attention is not not explanation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 11–20, 2019.

[23] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems*, pages 9505–9515, 2018.

[24] 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. *International Conference on Learning Representations*, 2(5):6, 2017.

[25] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *International Conference on Machine Learning*, pages 2654–2663, 2018.

[26] Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel. InfoGAN: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in neural information processing systems*, pages 2172–2180, 2016.

[27] Tejas D Kulkarni, William F Whitney, Pushmeet Kohli, and Josh Tenenbaum. Deep convolutional inverse graphics network. In *Advances in neural information processing systems*, pages 2539–2547, 2015.

[28] 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, 2019.

[29] Shalmali Joshi, Oluwasanmi Koyejo, Warut Vijitbenjaronk, Been Kim, and Joydeep Ghosh. Towards realistic individual recourse and actionable explanations in black-box decision making systems. *arXiv preprint arXiv:1907.09615*, 2019.

[30] Emily Denton, Ben Hutchinson, Margaret Mitchell, and Timnit Gebru. Detecting bias with generative counterfactual face attribute augmentation, 2019.

[31] Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without opening the black box: Automated decisions and the gdpr. *Harv. JL & Tech.*, 31:841, 2017.

[32] Maximilian Augustin, Alexander Meinke, and Matthias Hein. Adversarial robustness on in-and out-distribution improves explainability. *arXiv preprint arXiv:2003.09461*, 2020.

[33] Sarah R Beck, Kevin J Riggs, and Sarah L Gorniak. Relating developments in children's counterfactual thinking and executive functions. *Thinking & reasoning*, 15(4):337–354, 2009.

[34] Daphna Buchsbaum, Sophie Bridgers, Deena Skolnick Weisberg, and Alison Gopnik. The power of possibility: Causal learning, counterfactual reasoning, and pretend play. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 367(1599):2202–2212, 2012.

[35] Deena S Weisberg and Alison Gopnik. Pretense, counterfactuals, and bayesian causal models: Why what is not real really matters. *Cognitive science*, 37(7):1368–1381, 2013.

[36] Jeya Vikranth Jeyakumar, Joseph Noor, Yu-Hsi Cheng, Luis Garcia, and Mani Srivastava. How can i explain this to you? an empirical study of deep neural network explanation methods. *Advances in Neural Information Processing Systems*, 33, 2020.

[37] Julia K Winkler, Christine Fink, Ferdinand Toberer, Alexander Enk, Teresa Deinlein, Rainer Hofmann-Wellenhof, Luc Thomas, Aimilios Lallas, Andreas Blum, Wilhelm Stolz, et al. Association between surgical skin markings in dermoscopic images and diagnostic performance of a deep learning convolutional neural network for melanoma recognition. *JAMA dermatology*, 155(10):1135–1141, 2019.

[38] Michael Downs, Jonathan L Chu, Yaniv Yacoby, Finale Doshi-Velez, and Weiwei Pan. Cruds: Counterfactual recourse using disentangled subspaces. *ICML Workshop on Human Interpretability in Machine Learning*, 2020.

[39] Ramaravind K Mothilal, Amit Sharma, and Chenhao Tan. Explaining machine learning classifiers through diverse counterfactual explanations. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 607–617, 2020.

[40] David Alvarez Melis and Tommi Jaakkola. Towards robust interpretability with self-explaining neural networks. In *Advances in Neural Information Processing Systems*, pages 7775–7784, 2018.

[41] Sumedha Singla, Brian Pollack, Stephen Wallace, and Kayhan Batmanghelich. Explaining the black-box smoothly-a counterfactual approach. *arXiv preprint arXiv:2101.04230*, 2021.

[42] Gregory Plumb, Maruan Al-Shedivat, Ángel Alexander Cabrera, Adam Perer, Eric Xing, and Ameet Talwalkar. Regularizing black-box models for improved interpretability. *Advances in Neural Information Processing Systems*, 33, 2020.

[43] Amirata Ghorbani, Abubakar Abid, and James Zou. Interpretation of neural networks is fragile. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3681–3688, 2019.

[44] Chris Burgess and Hyunjik Kim. 3D shapes dataset, 2018. URL https://github.com/deepmind/3dshapes-dataset.

[45] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *IEEE international conference on computer vision*, pages 3730–3738, 2015.

[46] Leon Sixt, Maximilian Granz, and Tim Landgraf. When explanations lie: Why modified bp attribution fails. *arXiv preprint arXiv:1912.09818*, 2019.

[47] Yash Goyal, Uri Shalit, and Been Kim. Explaining classifiers with causal concept effect (CaCE). *arXiv preprint arXiv:1907.07165*, 2019.

[48] Runjin Chen, Hao Chen, Jie Ren, Ge Huang, and Quanshi Zhang. Explaining neural networks semantically and quantitatively. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9187–9196, 2019.

[49] Jose Oramas, Kaili Wang, and Tinne Tuytelaars. Visual explanation by interpretation: Improving visual feedback capabilities of deep neural networks. In *International Conference on Learning Representations*, 2018.

[50] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

[51] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.

[52] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier GANs. In *International Conference on Machine Learning*, volume 70, pages 2642–2651, 2017.

[53] Takeru Miyato and Masanori Koyama. cGANs with projection discriminator. In *International Conference on Learning Representations*, 2018. URL https://openreview.net/forum?id=ByS1VpgRZ.

[54] Rui Shu, Yining Chen, Abhishek Kumar, Stefano Ermon, and Ben Poole. Weakly supervised disentanglement with guarantees. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=HJgSwyBKvr.

[55] Francesco Locatello, Michael Tschannen, Stefan Bauer, Gunnar Rätsch, Bernhard Schölkopf, and Olivier Bachem. Disentangling factors of variations using few labels. In *International Conference on Learning Representations*, 2019.

[56] Jack Klys, Jake Snell, and Richard Zemel. Learning latent subspaces in variational autoencoders. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 6445–6455, 2018.

[57] Yash Goyal, Ziyan Wu, Jan Ernst, Dhruv Batra, Devi Parikh, and Stefan Lee. Counterfactual visual explanations. In *International Conference on Machine Learning*, pages 2376–2384, 2019.

[58] Amit Dhurandhar, Pin-Yu Chen, Ronny Luss, Chun-Chen Tu, Paishun Ting, Karthikeyan Shanmugam, and Payel Das. Explanations based on the missing: Towards contrastive explanations with pertinent negatives. In *Advances in neural information processing systems*, pages 592–603, 2018.

[59] Christopher P Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. Understanding disentangling in beta-VAE. *Advances in neural information processing systems: Workshop on Learning Disentangled Representations*, 2017.

[60] Tian Qi Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disentanglement in variational autoencoders. In *Advances in neural information processing systems*, pages 2610–2620, 2018.

[61] Abhishek Kumar, Prasanna Sattigeri, and Avinash Balakrishnan. Variational inference of disentangled latent concepts from unlabeled observations. In *International Conference on Learning Representations*, 2018. URL `https://openreview.net/forum?id=H1kG7GZAW`.

[62] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 6309–6318, 2017.

[63] Ali Razavi, Aaron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL `https://proceedings.neurips.cc/paper/2019/file/5f8e2fa1718d1bbcadf1cd9c7a54fb8c-Paper.pdf`.

[64] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2536–2544, 2016.

[65] JH Rick Chang, Chun-Liang Li, Barnabas Poczos, BVK Vijaya Kumar, and Aswin C Sankaranarayanan. One network to solve them all–solving linear inverse problems using deep projection models. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5888–5897, 2017.

[66] Mihaela Rosca, Balaji Lakshminarayanan, David Warde-Farley, and Shakir Mohamed. Variational approaches for auto-encoding generative adversarial networks. *arXiv preprint arXiv:1706.04987*, 2017.

[67] David Berthelot, Thomas Schumm, and Luke Metz. Began: Boundary equilibrium generative adversarial networks. *arXiv preprint arXiv:1703.10717*, 2017.

[68] Anh Nguyen, Jeff Clune, Yoshua Bengio, Alexey Dosovitskiy, and Jason Yosinski. Plug & play generative networks: Conditional iterative generation of images in latent space. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4467–4477, 2017.

[69] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.

[70] Junbo Zhao, Michael Mathieu, and Yann LeCun. Energy-based generative adversarial network. *arXiv preprint arXiv:1609.03126*, 2016.

[71] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.

[72] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. Improved training of wasserstein gans. *arXiv preprint arXiv:1704.00028*, 2017.

[73] Chun-Liang Li, Wei-Cheng Chang, Yu Cheng, Yiming Yang, and Barnabás Póczos. Mmd gan: Towards deeper understanding of moment matching network. *arXiv preprint arXiv:1705.08584*, 2017.

[74] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.

[75] Zinan Lin, Kiran Thekumparampil, Giulia Fanti, and Sewoong Oh. InfoGAN-CR and ModelCentrality: Self-supervised model training and selection for disentangling GANs. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 6127–6139. PMLR, 13–18 Jul 2020. URL `http://proceedings.mlr.press/v119/lin20e.html`.

[76] Hensin Tsao, Jeannette M Olazagasti, Kelly M Cordoro, Jerry D Brewer, Susan C Taylor, Jeremy S Bordeaux, Mary-Margaret Chren, Arthur J Sober, Connie Tegeler, Reva Bhushan, et al. Early detection of melanoma: reviewing the abcdes. *Journal of the American Academy of Dermatology*, 72(4):717–723, 2015.

[77] Darrell S Rigel, Robert J Friedman, Alfred W Kopf, and David Polsky. Abcde—an evolving concept in the early detection of melanoma. *Archives of dermatology*, 141(8):1032–1034, 2005.

[78] Thomas B Fitzpatrick. The validity and practicality of sun-reactive skin types i through vi. *Archives of dermatology*, 124(6):869–871, 1988.

[79] Yi Sun, Yuheng Chen, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation by joint identification-verification. In *Proceedings of the 27th International Conference on Neural Information Processing Systems-Volume 2*, pages 1988–1996, 2014.

[80] Suman Ravuri and Oriol Vinyals. Classification accuracy score for conditional generative models. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.

[81] Yanai Elazar and Yoav Goldberg. Adversarial removal of demographic attributes from text data. In *Conference on Empirical Methods in Natural Language Processing*, pages 11–21, 2018.

[82] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in neural information processing systems*, pages 6626–6637, 2017.

[83] Tim Salimans, Ian J Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training GANs. In *NIPS*, 2016.

[84] Jenna Wiens, Suchi Saria, Mark Sendak, Marzyeh Ghassemi, Vincent X Liu, Finale Doshi-Velez, Kenneth Jung, Katherine Heller, David Kale, Mohammed Saeed, et al. Do no harm: A roadmap for responsible machine learning for health care. *Nature medicine*, pages 1–4, 2019.

[85] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. Model cards for model reporting. In *Proceedings of the conference on fairness, accountability, and transparency*, pages 220–229, 2019.

[86] Marilyn S Sommers, Jamison D Fargo, Yadira Regueira, Kathleen M Brown, Barbara L Beacham, Angela R Perfetti, Janine S Everett, and David J Margolis. Are the fitzpatrick skin phototypes valid for cancer risk assessment in a racially and ethnically diverse sample of women? *Ethnicity & disease*, 29(3):505, 2019.

[87] Latrice C Pichon, Hope Landrine, Irma Corral, Yongping Hao, Joni A Mayer, and Katherine D Hoerster. Measuring skin cancer risk in african americans: is the fitzpatrick skin type classification scale culturally sensitive. *Ethn Dis*, 20(2):174–179, 2010.

[88] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International Conference on Learning Representations*, 2015.

# A Appendix

## A.1 Baselines

### A.1.1 Baseline 1: multi-modal explainability through VAE-based disentanglement

Disentanglement approaches have demonstrated practical success in learning representations that correspond to factors for variation in data [54], though some gaps between theory and practice remain [28]. However, the extent to which these techniques can aid post hoc explainability in conjunction with an external model is not well understood. Thus, we consider a set of baseline approaches based

on VAEs explicitly designed for disentanglement: $\beta$-VAE [24], Annealed-VAE [59], $\beta$-TCVAE [60], and DIPVAE [61]. We extend each of them to incorporate the classifier's signal during their training processes for a fair comparison with DISSECT. Intuitively speaking, this encourages the generative model to learn latent dimensions that could influence the classifier, i.e., learning *Influential* CTs. Note that it is necessary to also cover data generative factors that are independent of the external classifier. This means that, along with having a good quality reconstruction and high likelihood, we need to promote sparsity in sensitivity of the external classifier to the latent dimensions of the generative model.

More formally, consider a vanilla VAE that has an encoder $e_\theta$ with parameters $\theta$, a decoder $d_\phi$ with parameters $\phi$, and the $M$-dimensional latent code $z$ with prior distribution $p(z)$. Recall that $x$ denotes the input sample. The objective of a VAE is to minimize the loss:

$$\mathcal{L}_{\theta,\phi}^{\text{vanilla VAE}} = -\mathbb{E}_{z \sim e_\theta(z|x)}[\log d_\phi(x|z)] + \mathbb{KL}(e_\theta(z|x)||p(z)).$$

We introduce an additional loss term for incorporating the black-box classifier's signal: $1/K \sum_{k=1}^K \partial f(x)/\partial z_k$. We impose this only for the first $K$ dimensions in the latent space[4], in other words, the number of desired CTs, $K \leq M$. Minimizing this term provides $\text{CT}_k$s, $k \in \{1, 2, \ldots, K\}$, with negative $\partial f(x)/\partial z_k$, with a high $|\partial f(x)/\partial z_k|$. The final loss is:

$$\mathcal{L}_{\theta,\phi} = \mathcal{L}_{\theta,\phi}^{\text{vanilla VAE}} + \lambda * \frac{\sum_{k=1}^K \partial f(x)/\partial z_k}{K},$$

where $\lambda$ is a hyper-parameter. We apply this modification to the four aforementioned VAE-based approaches and refer to them with a -mod postfix, e.g., $\beta$-VAE-mod.[5]

**Development process.** To promote discovering *Important* CTs, we introduced $\mathcal{L}_{\text{aux}} = \frac{\sum_{d=1}^K \partial f(x)/\partial z_d}{K}$, which incorporated the directional derivative of $f$ with respect to the latent dimensions of interest into the loss function of VAE. Despite experimentation with many variants of $\mathcal{L}_{\text{aux}}$, we observed two common themes.

First, a monotonic increase of $f(\bar{x})$ through traversing one latent dimension and keeping the rest static was hardly achieved. Second, while the purpose of $\mathcal{L}_{\text{aux}}$ was to promote exerting *Importance* only in the first $K$ dimensions of the latent space, $\partial f(x)/\partial z_d$ for $d \in \{K+1, K+2, \cdots, M\}$ were impacted similarly. Having strongly correlated dimensions is a failure in achieving the very goal of disentanglement approaches. Table 4 summarizes a subset of the variants of $\mathcal{L}_{\text{aux}}$ studied.

Table 4: Summary of a subset of $\mathcal{L}_{\text{aux}}$ iterations. The development goal is to make the first $K$ dimensions of the latent space *Important*. In some iterations, we encouraged the remaining $M - K$ dimensions not to be *Important* to reduce potential correlation across latent dimensions.

| | $\mathcal{L}_{\text{aux}}$ |
|---|---|
| 1 | $\frac{\sum_{k=1}^K \partial f(x)/\partial z_k}{K}$ |
| 2 | $\frac{\sum_{k=1}^K \partial f(x)/\partial z_k}{K} + \frac{\sum_{d=K+1}^M |\partial f(x)/\partial z_d|}{M-K}$ |
| 3 | $\frac{\sum_{k=1}^K \partial f(x)/\partial z_k}{K} + \frac{\sum_{d=K+1}^M [\partial f(x)/\partial z_d]^2}{M-K}$ |
| 4 | $\frac{\sum_{d=K+1}^M |\partial f(x)/\partial z_d|}{M-K}$ |
| 5 | $\frac{\sum_{d=K+1}^M [\partial f(x)/\partial z_d]^2}{M-K}$ |
| 6 | $\frac{\sum_{k,d} |\partial f(x)/\partial z_d|/|\partial f(x)/\partial z_k|}{K*(M-K)}$ where $k \in \{1, 2, \cdots, K\}, d \in \{K+1, K+2, \cdots, M\}$ |
| 7 | $\frac{\sum_{k,d} log(|\partial f(x)/\partial z_d|/|\partial f(x)/\partial z_k|)}{K*(M-K)}$ where $k \in \{1, 2, \cdots, K\}, d \in \{K+1, K+2, \cdots, M\}$ |

---

[4]Without loss of generality, the additional term can be applied to the first $K$ dimensions, and there is no need to consider $\binom{K}{M}$ potential selections.

[5]We build these baselines on top of the open-sourced implementations provided in `https://github.com/google-research/disentanglement_lib` under Apache License 2.0.

### A.1.2 Baseline 2: multi-modal explainability through conditional subspace VAE

Another relevant area of work is conditional generation. In particular, Conditional subspace VAE (CSVAE) is a method aiming to solve unsupervised learning of features associated with a specific label using a low-dimensional latent subspace that can be independently manipulated [56]. CSVAE partitions the latent space into two parts: $w$ learns representations correlated with the label, and $z$ covers the remaining characteristics for data generation. An assumption of independence between $z$ and $w$ is made. To explicitly enforce independence in the learned model, we minimize the mutual information between $Y$ and $Z$. CSVAE has proven successful in providing counterfactual scenarios to reverse unfavorable decisions of an algorithm, also known as algorithmic recourse [38]. To adjust CSVAE to explain the decision-making of an external classifier $f$, we treat the predictions of the classifier as the label of interest.

More formally, the generative model can be summarized as:

$$w|y \sim N(\mu_y, \sigma_y^2.I), y \sim Bern(p),$$
$$x|w, z \sim N(d_{\phi_\mu}(w, z), \sigma_\epsilon^2.I), z \sim N(0, \sigma_z^2.I)$$

Conducting inference leads to the following objective function:

$$M_1 = \mathbb{E}_{D(x,y)}[-\mathbb{E}_{q_\phi(z,w|x,y)}[\log p_\theta(x|w,z)] + \mathbb{KL}(q_\phi(w|x,y)||p_\gamma(w|y)) + \mathbb{KL}(q_\phi(z|x,y)||p(z)) - \log p(y)]$$

$$M_2 = \mathbb{E}_{q_\phi(z|x)} D(x)[\int_Y q_\delta(y|z) \log q_\delta(y|z) \, dy]$$

$$M_3 = \mathbb{E}_{q(z|x)D(x,y)}[q_\delta(y|z)]$$
$$\min_{\theta,\phi,\gamma} \beta_1 M_1 + \beta_2 M_2; \quad \max_\delta \beta_3 M_3$$

### A.1.3 Baseline 3: multi-modal explainability through progressive exaggeration

Explanation by Progressive Exaggeration (**EPE**) [9] is a recent successful generative approach that learns to generate one series of counterfactual and realistic samples that change the prediction of $f$, given data and the classifier's signal. It is particularly relevant to our work as it explicitly optimizes *Influence* and *Realism*. EPE is a type of Generative Adversarial Network (GAN) [51] consisting of a discriminator ($D$) and a generator ($G$) that is based on Projection GAN [53]. It incorporates the amount of desired perturbation $\alpha$ on the outcome of $f$ as:

$$\mathcal{L}_{\text{cGAN}}(D) = -\mathbb{E}_{x \sim p_{data}}[\min(0, -1 + D(x, 0))] - \mathbb{E}_{x \sim p_{data}}[\min(0, -1 - D(G(x, \alpha), \alpha))] \quad (1)$$

$$\mathcal{L}_{\text{cGAN}}(G) = -\mathbb{E}_{x \sim p_{data}}[D(G(x, \alpha), \alpha)] \quad (2)$$

A Kullback–Leibler divergence (KL) term in the objective function between the desired perturbation ($\alpha$) and the achieved one ($f(G(x, \alpha))$) promotes Importance [9]:

$$\mathcal{L}_f(D, G) = r(D, G(x, \alpha)) + \mathbb{KL}(\alpha|f(G(x, \alpha))), \quad (3)$$

where the first term is the likelihood ratio defined in projection GAN [51], and [9] uses an ordinal-regression parameterization of it.

A reconstruction loss and a cycle loss promote self-consistency in the model, meaning that applying a reverse perturbation or no perturbation should reconstruct the query sample:

$$\mathcal{L}_{\text{rec}}(G) = ||x - G(x, f(x))||_1 \quad (4)$$

$$\mathcal{L}_{\text{cyc}}(G) = ||x - G(G(x, \alpha), f(x))||_1. \quad (5)$$

Thus, the overall objective function of EPE is the following:

$$\min_G \max_D \lambda_{\text{cGAN}} \mathcal{L}_{\text{cGAN}}(D, G) + \lambda_f \mathcal{L}_f(D, G) + \lambda_{\text{rec}} \mathcal{L}_{\text{rec}}(G) + \lambda_{\text{rec}} \mathcal{L}_{\text{cyc}}(G), \quad (6)$$

where $\lambda_{\text{cGAN}}$, $\lambda_f$, and $\lambda_{\text{rec}}$ are the hyper-parameters.

Note that EPE only finds one pathway to switch the classifier's outcome. We argue that classifiers learned from challenging and realistic datasets will have complex reasoning pathways that could enhance model explainability if revealed. Decomposing this complexity is needed to make reasoning comprehensible for humans. We thus create a more powerful baseline, an EPE-variant, **EPE-mod**. EPE-mod learns multiple pathways by making the generator conditional on another variable: the CT dimension. More formally, EPE-mod updates $G(\cdot, \cdot)$ to $G(\cdot, \cdot, k)$ in Eq. (1)-(5), while Eq. (6) remains unchanged. We compare DISSECT to both EPE and EPE-mod.
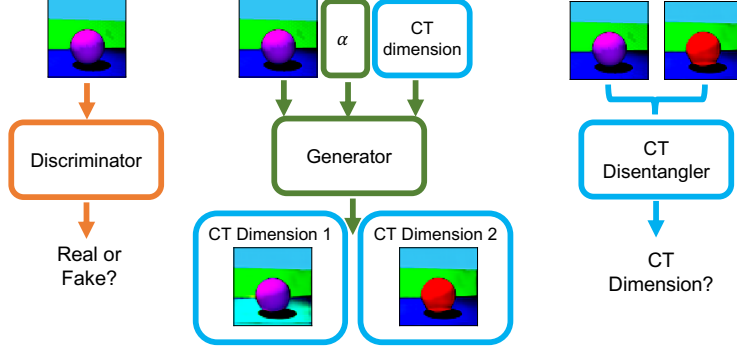
Figure 6: Simplified illustration of DISSECT. Orange, Green, and Blue show elements related to the discriminator, generator, and CT disentangler, respectively. DISSECT builds on top of [9] and has the Orange and Green components in common with it.

## A.2 DISSECT details

We build our proposed method on EPE-mod and further promote distinctness across CTs by adding a disentangler network, $R$. The disentangler is a classifier with K classes. Given a pair of $< x, x' >$ images, $R$ tries to predict which $CT_{k;k \in \{1,...,K\}}$ has perturbed query $x$ to produce $x'$. Note that $R$ can return close to 0 probability for all classes if $x'$ is just a reconstruction of $x$, indicating no tweaked dimensions. The disentangler also penalizes the generator if any CTs use similar pathways to cross the decision boundary. See the appendix for schematics of our method.

To formalize this, let:

$$\hat{x}_k = G(x, f(x), k)$$
$$\bar{x}_k = G(x, \alpha, k)$$
$$\tilde{x}_k = G(\bar{x}_k, f(x), k).$$

Note that $\hat{x}_k$ and $\tilde{x}_k$ are reconstructions of $x$ while $\bar{x}_k$ is perturbed to change the classifier output from $f(x)$ to $\alpha$. Therefore, $x, \bar{x}_k$ and $\tilde{x}_k$ form a cycle, and $k$ represents $CT_k$. $R(.,.)$ is the predicted probabilities of the perturbed concept given a pair of examples, which is a vector of size $K$, where each element is a value in $[0, 1]$. We define the following cross entropy loss that is a function of both $R$ and $G$:

$$\mathcal{L}_r(G, R) = CE(e_k, R(x, \bar{x}_k)) + CE(e_k, R(\bar{x}_k, \tilde{x}_k))$$
$$= -\mathbb{E}_{x \sim p_{data}} \sum_{k=1}^{K} [e_k log R(x, \bar{x}_k) + e_k log R(\bar{x}_k, \tilde{x}_k)] \tag{7}$$

Here, $e_k$ refers to a one-hot vector of size $K$ where the $k$-th element is one and the remaining elements are zero. This term enforces $R$ to identify no change when receiving reconstructions of the same image as input and utilizes the cycle and promotes determining the correct dimension when a non-zero change has happened, either increasing or decreasing the outcome of $f$. In summary, the overall objective function of our method is:

$$\min_{G,R} \max_{D} [\lambda_{cGAN} \mathcal{L}_{cGAN}(D, G) + \lambda_f \mathcal{L}_f(D, G) + \lambda_{rec} \mathcal{L}_{rec}(G) + \lambda_{rec} \mathcal{L}_{cyc}(G)$$
$$+ \lambda_r \mathcal{L}_r(G, R)] \tag{8}$$

For this adversarial min-max optimization, we use the Adam optimizer [88]. See Figure 6 for a for a simplified visualization of DISSECT or Figure 7 for a detailed version. We open-source our implementation at `https://github.com/asmadotgh/dissect` under MIT license, which builds on top of the open-source implementation of EPE [9] [6].

---

[6] `https://github.com/batmanlab/Explanation_by_Progressive_Exaggeration` available under MIT License.

Figure 7: Illustration of DISSECT. Orange, Green, and Blue show elements related to the discriminator, generator, and CT disentangler, respectively. DISSECT builds on top of [9] and has the Orange and Green components in common with it.

### A.3 Evaluation metrics details

The VAE-based baselines support continuous values for latent dimensions $z_k$. Also, we can directly sample latent code values and produce $CT_k$ by keeping $z_j$ ($j \neq k$) constant and monotonically increasing $z_k$ values. However, to calculate the evaluation metrics comparably to EPE, EPE-mod, and DISSECT, we do the following: We encode each query sample using the probabilistic encoder. We set $z_j = \mu_j$, $j \neq k$ where $\mu_j$ is the mean of the fitted Gaussian distribution for $z_j$. For dimension $k$, we produce $N + 1$ linearly spaced values between $\mu_p \pm 2 * \sigma_p$, where $\mu_p$ and $\sigma_p$ are the mean and standard deviation of the prior normal distribution, in our case 0.0 and 1.0 respectively. Note that these different values for $z_k$ map out to $\alpha$, $\alpha \in \{0, \frac{1}{N}, \cdots, 1\}$ in EPE, EPE-mod, and DISSECT models. After this step, calculating all the metrics related to *Importance*, *Realism*, and *Distinctness* is identical across all the models.

### A.4 Experiment setup and hyper-parameter tuning details

Experiments were conducted on an internal compute cluster at MIT Media Lab. Training and evaluation of all models across the three datasets have approximately taken 1000 hours on a combination of Nvidia GPUs including GTX TITAN X, GTX 1080 Ti, RTX 2080 Rev, and Quadro K5200.

We seeded the model's parameters from [9] based on the reported values in their accompanying open-sourced repository[7]. We used the same parameters for 3D Shapes, except for the number of bins, $N$, used for ordinal regression transformation of the classifier's posterior probability. The largest number of bins that resulted in non-zero samples per bin, 3, was selected. We kept all the parameters shared between EPE, EPE-mod, and DISSECT the same.

Given the experiments' design, we fixed the number of dimensions $K$ in DISSECT and EPE-mod to 2. We experimented with a few values for $\lambda_r$, 1, 10, 20, 50. Based on manual inspection after 30k training batches, $\lambda_r$ was selected. Factors considered for selection included inspecting the perceived quality of generated samples and the learning curves of $\mathcal{L}_{cGAN}(D)$, $\mathcal{L}_{cGAN}(G)$, $\mathcal{L}_{cyc}(G)$, $\mathcal{L}_{rec}(G)$, and $\mathcal{L}_r(G, R)$.

---

[7]https://github.com/batmanlab/Explanation_by_Progressive_Exaggeration available under MIT License.

Table 5: Summary of hyper-parameter values. Discriminator optimization happens once every $D$ steps. Similarly, generator optimization happens once every $G$ steps. $\lambda_r$ is specific to DISSECT, and $K$ is specific to EPE-mod and DISSECT. All the remaining parameters are shared across EPE, EPE-mod, and DISSECT. Note that samples used for evaluation are not included in the training process.

| | Preprocessing | | Training | | | | | | | | | Evaluation Metrics | | | |
| | $N$ | max samples per bin | $\lambda_{cGAN}$ | $\lambda_{rec}$ | $\lambda_f$ | $D$ steps | $G$ steps | batch size | epochs | $K$ | $\lambda_r$ | max # samples | batch size | epochs | hold-out test ratio |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3D Shapes | 3 | 5,000 | 1 | 100 | 1 | 1 | 5 | 32 | 300 | 2 | 10 | 10,000 | 32 | 10 | 0.25 |
| SynthDerm | 2 | 1,350 | 2 | 100 | 1 | 5 | 1 | 32 | 300 | 5 | 2 | 10,000 | 8 | 10 | 0.25 |
| CelebA | 10 | 5,000 | 1 | 100 | 1 | 1 | 5 | 32 | 300 | 2 | 10 | 10,000 | 32 | 10 | 0.25 |



Figure 8: Qualitative results on 3D Shapes when flipping classification outcome from "False" to "True." We observe that EPE-mod converges to finding the same concept, despite having the ability to express multiple pathways to switch the classifier outcome. However, DISSECT can discover the two Distinct ground-truth concepts: $CT_1$ flips the floor color to cyan, and $CT_2$ converts the shape color to red.
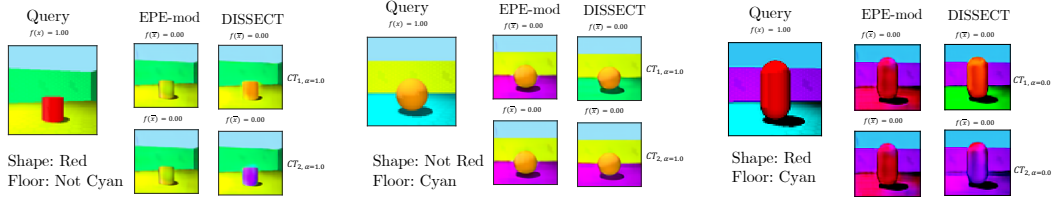


Figure 9: Qualitative results on 3D Shapes when flipping classification outcome from "True" to "False." We observe that EPE-mod converges to finding the same concept, despite having the ability to express multiple pathways to switch the classifier outcome. However, DISSECT is capable of discovering Distinct paths to do so. **Left**: When the input query has a red shape, but the floor color is not cyan, $CT_1$ flips the shape color to orange and $CT_2$ flips it to violet. **Middle**: When the input query has a cyan floor, but the shape color is not red, $CT_1$ flips the floor color to lime, and $CT_2$ converts it to magenta. **Right**: When the input query has a red shape and cyan floor, $CT_1$ changes the shape color to dark orange and floor color to lime, and $CT_2$ flips the shape color to violet and floor color to magenta.

For evaluation, we used a hold-out set including 10K samples. For post hoc evaluation classifiers predicting *Distinctness* and *Realism*, 75% of the samples were used for training, and the results were reported on the remaining 25%. See Table 5 for the summary of the hyper-parameter values.
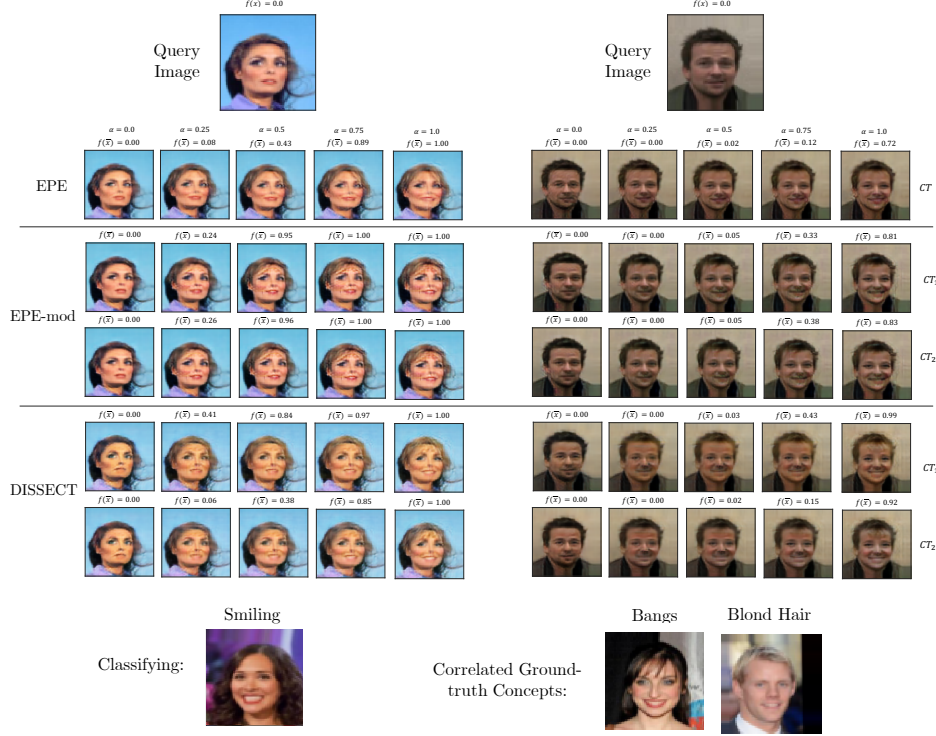
Figure 10: Qualitative results on `CelebA`. A biased classifier has been trained to predict smile probability, where the training dataset has been sub-sampled such that smiling co-occurs only with "bangs" and "blond hair" attributes. EPE does not support multiple CTs. We observe that EPE-mod converges to finding the same concept, despite having the ability to express several pathways to change $f(\bar{x})$ through $CT_1$ and $CT_2$. However, DISSECT can discover Distinct routes: $CT_1$ mainly changes hair color to blond, and $CT_2$ does not alter hair color but focuses more on hairstyle and tries to add bangs. Thus it identifies two otherwise hidden biases.

## A.5 Additional qualitative results for case study I

Recall that considering `3D Shapes`, we define an image as "colored correctly" if the shape hue is red *or* the floor hue is cyan. Given a not "colored correctly" query, we recover a CT related to the shape color and another CT associated with the floor color–two different pathways leading to switching the classifier outcome for that sample. See Figure 8 for additional qualitative examples where classification outcome is flipped from False to True.

However, these two ground-truth concepts do not directly apply to switching the classifier outcome from True to False in this scenario. For example, if an image has a red shape *and* a cyan floor, both colors need to be changed to switch the classification outcome. As shown in Figure 9, we still observe that applying DISSECT to such cases results in two discovered CTs that change different combinations of colors while EPE-mod converges to the same CT.

## A.6 Additional qualitative results for case study III

Recall the biased `CelebA` experiment where smiling correlates with "blond hair" and "bangs" attributes. Figure 10 shows additional qualitative samples, suggesting that DISSECT can recover and separate the aforementioned concepts, which other techniques fail to do.

## A.7 Additional influence metrics

Following [9], we conduct a more granular analysis to investigate if DISSECT works similarly across different queries, e.g., when flipping classification outcome from "True" to "False" or the other way around. We plot $\alpha$ vs. $f(\bar{x})$ for query samples with $f(x) < 0.5$ and $f(x) \geq 0.5$ separately.
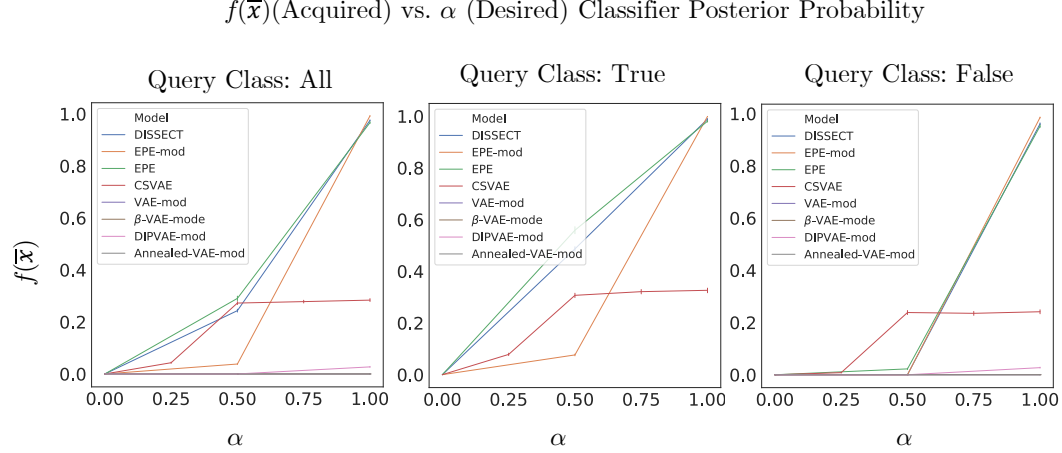
### A.7.1 Additional quantitative results for case study I

$f(\overline{x})$(Acquired) vs. $\alpha$ (Desired) Classifier Posterior Probability



Figure 11: Acquired vs. desired classifier posterior probability for generated samples that constitute a CT on `3D Shapes` over 10K queries in total. The ideal would be a line of slope one. Error bars represent 95% confidence intervals. We observe that DISSECT performs similarly to EPE that has been particularly geared toward exhibiting *Influence*, and its extension, EPE-mod. VAE based methods perform poorly in terms of *Influence*. CSVAE performs significantly better than other VAE baselines but still works much worse than EPE, EPE-mod, and DISSECT. There is a significant correlation between acquired and desired posterior probabilities of generated samples for DISSECT (r=0.82, p<.0001), EPE-mod (r=0.87, p<.0001), EPE (r=0.81, p<.0001), and CSVAE (r=0.32, p<.0001). In other VAE baselines, there is very low or no correlation between acquired and desired probabilities: DIPVAE (r=0.14, p<.0001), VAE (r=0.07, p<.0001), $\beta$-VAE-mode (r=-0.01, p>.1) and Annealed-VAE-mod (r=-0.01, p>.1).

Figure 11 depicts more details regarding CTs' *Importance* across different groups of samples for `3D Shapes` experiments.

### A.7.2 Additional quantitative results for case study II

Figure 12 provides more granular information about *Importance* scores that further confirms our qualitative results for `SynthDerm` experiments.

### A.7.3 Additional quantitative results for case study III

Figure 13 provides further details regarding *Importance* scores on a more granular scale for `celebA` dataset.

$f(\overline{x})$(Acquired) vs. $\alpha$ (Desired) Classifier Posterior Probability
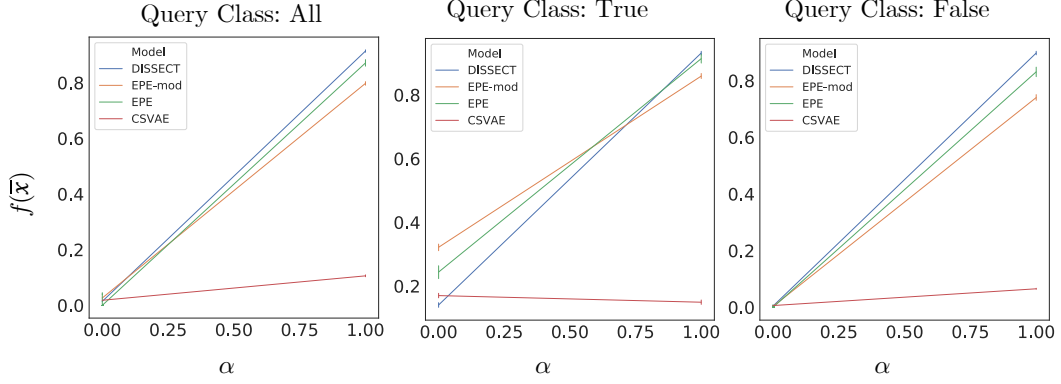
Figure 12: Acquired vs. desired classifier posterior probability for generated samples that constitute a CT on `SynthDerm` over 10K queries in total. The ideal would be a line of slope one. Error bars represent 95% confidence intervals. We observe that DISSECT performs similarly to EPE that has been particularly geared toward exhibiting *Influence*, and it potentially outperforms EPE-mod. Although CSVAE produces examples with acquired posterior probabilities correlated with the desired values (r=0.25, p<.0001), it performs significantly worse than EPE (r=0.87, p<.0001), EPE-mod (r=0.81, p<.0001), and DISSECT (r=0.92, p<.0001).



$f(\overline{x})$(Acquired) vs. $\alpha$ (Desired) Classifier Posterior Probability
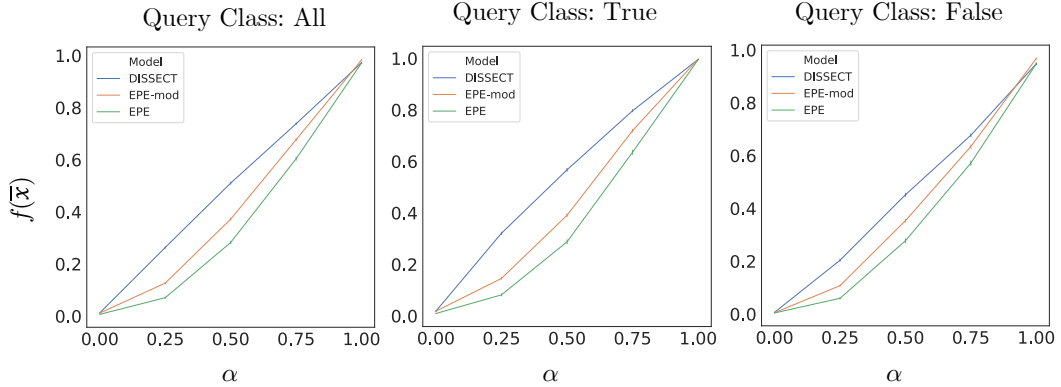
Figure 13: Acquired vs. desired classifier posterior probability for generated samples that constitute a CT on `CelebA` over 10K queries in total. The ideal would be a line of slope one. Error bars represent 95% confidence intervals. The results suggest that DISSECT performs on par with the three strongest baselines in terms of *Importance*. Acquired and desired probabilities of generated samples are significantly correlated for DISSECT (r=0.84, p<.0001), EPE-mod (r=0.86, p<.0001), and EPE (r=0.85, p<.0001).