

---

# Uncovering Conceptual Blindspots in Generative Image Models Using Sparse Autoencoders

---

**Matyas Bohacek<sup>1\*</sup>**   **Thomas Fel<sup>2\*</sup>**   **Maneesh Agrawala<sup>1</sup>**   **Ekdeep Singh Lubana<sup>3</sup>**

<sup>1</sup>Department of Computer Science, Stanford University

<sup>2</sup>Kempner Institute, Harvard University

<sup>3</sup>CBS-NTT Program in Physics of Intelligence, Harvard University

## Abstract

Despite their impressive performance, generative image models trained on large-scale datasets frequently fail to produce images with seemingly simple concepts—e.g., human hands or objects appearing in groups of four—that are reasonably expected to appear in the training data. These failure modes have largely been documented anecdotally, leaving open the question of whether they reflect idiosyncratic anomalies or more structural limitations of these models. To address this, we introduce a systematic approach for identifying and characterizing “conceptual blindspots”—concepts present in the training data but absent or misrepresented in a model’s generations. Our method leverages sparse autoencoders (SAEs) to extract interpretable concept embeddings, enabling a quantitative comparison of concept prevalence between real and generated images. We train an archetypal SAE (RA-SAE) on DINOv2 features with 32,000 concepts—the largest such SAE to date—enabling fine-grained analysis of conceptual disparities. Applied to four popular generative models (Stable Diffusion 1.5/2.1, PixArt, and Kandinsky), our approach reveals specific suppressed blindspots (e.g., bird feeders, DVD discs, and whitespaces on documents) and exaggerated blindspots (e.g., wood background texture and palm trees). At the individual datapoint level, we further isolate memorization artifacts — instances where models reproduce highly specific visual templates seen during training. Overall, we propose a theoretically grounded framework for systematically identifying conceptual blindspots in generative models by assessing their conceptual fidelity with respect to the underlying data-generating process.

## 1 Introduction

Generative image models trained on large scale datasets have achieved unprecedented capabilities, allowing their use in applications both within the vision domain [1–9] and well beyond that [9–15]. Despite this success, several qualitative [16–18] and quantitative studies [19, 20] have shown that, at times, models can struggle to generate images with relatively simple concepts, e.g., human hands [21–24], objects appearing in groups of four [25], and negations or object relations [20, 26]. In fact, when prompted to generate images containing such concepts, models tend to produce outputs with related structures, but not precisely the ground truth concept—e.g., producing hands with six fingers. These failure modes, which we call “conceptual blindspots”<sup>2</sup>, can be unintuitive, since one may reasonably expect models have had enough exposure

---

\*Equal contribution. Email: maty@stanford.edu, tfel@g.harvard.edu, maneesh@cs.stanford.edu, ekdeplubana@fas.harvard.edu.

<sup>2</sup>We borrow the term “blindspots” from psychology literature [27], wherein it is used to refer to scenarios where an agent makes biased decisions despite exposure to observations that contradicts their decision’s rationale.



Figure 1:  $\delta(k)$  quantifies a model’s tendency to over- or under-generate a concept  $c_k$  compared to its natural-data frequency. We deem concepts with  $\delta(k) < 0.1$  as **suppressed conceptual blindspots** and concepts with  $\delta(k) > 0.9$  as **exaggerated conceptual blindspots**. The depicted images, generated by four popular generative image models, show examples of images with conceptual blindspots, as well as aligned concepts.

to demonstrations accurately detailing such concepts. This raises the question whether such failures reflect intriguing quirks of certain specific concepts, or whether they are demonstration of a more systematic phenomenon under which, for a broad spectrum of concepts, models fail to or are overly likely to produce images containing them.

Answering this question remains infeasible with existing approaches for evaluating generative image models [28, 29]. Specifically, existing approaches generally rely on coarse-grained measures that are meant to assess image realism, e.g., FID [30], and hence do not capture distributional failures. Methods like CLIPScore evaluate generation diversity or distribution-coverage statistics [31–33], hence offering partial insights to our question, but not at the granularity of fine-grained features or concepts [34, 35], making it difficult to identify conceptual blindspots. Finally, qualitative analyses for evaluating generative models, such as participant surveys [5, 17, 36–38] or open-ended exploration [39], can identify failures in models’ ability to capture certain concepts, but do not offer a scalable methodology that can be easily repeated across model families and used for their comparison.

**This work.** Motivated by the above, we argue identifying and analyzing conceptual blindspots in a generative image model requires designing a methodology that, in an automated and unsupervised manner, can elicit concept in the data distribution that have a mismatch between their odds of generation by the true data-generating process versus the trained model. Our contributions in this work are three-fold:

- **Formalizing Conceptual Blindspots in Generative Image Models.** We introduce a systematic framework for identifying and quantifying conceptual blindspots in generative image models compared to natural images (Section 3). This formalization moves beyond anecdotal or human-defined concept evaluations, offering a principled approach to understand the models’ limitations.
- **A Scalable, Unsupervised Approach for Identifying Conceptual Blindspots using Sparse Autoencoders.** We develop a structured methodology using sparse autoencoders (SAEs) to extract and compare concept distributions between natural and synthesized images (Section 4). To this end, we train and open-source an archetypal SAE (RA-SAE) on DINOv2 features with 32,000 concepts, the largest such RA-SAE to date.
- **Exploratory Tool and Analysis.** Our exploratory web tool enables both distribution- and datapoint-level analysis of blindspots across models (Sections 5.1–5.4). We apply our method to Kandinsky, PixArt, and Stable Diffusion (SD) 1.5/2.1 (Section 5). We identify specific instances of both suppressed conceptual blindspots (e.g., bird feeders, DVD discs, and whitespaces on documents) and exaggerated conceptual blindspots (e.g., wood background texture and palm trees), shown in Fig. 1.

The model, code, and web tool are available at <https://conceptual-blindspots.github.io/>.

## 2 Related Work

### 2.1 Explainability in Vision

Early work in explainable AI, including computer vision, focused on methods for attribution of influential input regions [40–42]. However, these methods offered limited semantic information about learned representations and often produced incorrect explanations [43–45]. To address these issues,

concept-based interpretability [46] emerged to identify semantically meaningful directions in neural networks, revealing not just where they look but what concepts and structures they employ [47–49].

Recent work demonstrates that popular concept-based interpretability methods—ACE [44], CRAFT [48], and Sparse Autoencoders (SAEs)[50, 51]—essentially address the same dictionary learning task under different constraints[52]. Out of these approaches, sparse autoencoders (SAEs) have emerged as particularly scalable for concept-based interpretability. While recent studies reveal some limitations of the original SAEs—including overly specific features [53], compositionality challenges [54], and limited intervention effects [55]—improved SAE versions have emerged, including archetypal SAE (RA-SAE) [56].

Beyond SAEs, other interpretability methods include prompt-based probing [57], attention map or activation visualizations [58, 59], and dataset-level statistics [32, 33] (e.g., diversity or distribution coverage metrics) offer only partial insights to answer these questions. Crucially, they lack granularity, focusing on full images or prompts instead of fine-grained features and concepts [34, 35]. Furthermore, they depend on subjective interpretation and do not distinguish between various failure models [60]. The existing methods and metrics are hence inadequate in systematically identifying feature- and concept-level weaknesses of generative image models [28].

## 2.2 Generative Image Models

Diffusion-based methods have become dominant across various modalities in generative vision modeling, including image [4, 61–64], video [65–72], and 3D [7, 73–75]. In the domain of image generation, this can be traced back to denoising diffusion probabilistic models (DDPMs) [76], which were later extended to non-Markov diffusion processes with denoising diffusion implicit models (DDIMs) [77].

The Stable Diffusion (SD) [9] model family made DDPMs highly accessible both in the research and open-source communities. The original SD was followed by several subsequent versions, including SD 2 [78], SD 3 [79], SD XL [80]. Many modifications and extensions of the SD architecture have emerged, enabling additional constraints for the diffusion process (e.g., style [81, 82], pose [83], and identity [84, 85]) as well as different input modalities, such as image-to-image generation. Different Latent Diffusion Models built on top of SD—including Kandinsky [86], PixArt [87], DeepFloyd [4, 88], and FLUX [89]—have also emerged.

## 2.3 Datasets for Generative Image Models

The recent success of generative vision models is largely attributed to the abundance of computational resources and large-scale internet datasets [90–92]. Specifically, LAION-5B [93] has played a key role in the training of open-source text-to-image models, including SD and its derivatives. This dataset, scraped from Common Crawl [94], contains over 5B image-caption pairs, 2.3B of which are in English. Other prominent datasets include COYO-700M [95] and Conceptual Captions [96], with 700M and 3M image-caption pairs, respectively.

As LAION-5B gained popularity, concerns grew over its biases [97–101]. Despite filtering attempts, harmful content persisted [97–99], including NSFW material [100] and hundreds of CSAM instances [101], prompting its temporary removal from official channels. The dataset also suffers from low-quality images [102] and internet-style captions (e.g., product descriptions) that misalign with how users prompt trained models [103].

## 2.4 Concept Discovery and Sparse Coding in Generative Image Models

Dictionary learning seeks to find sparse representations of input data, where each sample can be reconstructed using a linear combination of few dictionary atoms [104–106]. Built upon compressed sensing theory [107, 108], the field evolved from early vector quantization methods [109] to sophisticated approaches including Non-negative Matrix Factorization [110, 111], Sparse PCA [112], and K-SVD [113]. Recent advances include online methods [114], structured sparsity [115], and theoretical guarantees [116, 117], alongside growing connections to deep learning [118, 119].

Advances in sparse coding have also been leveraged to study the emergence of high-level concepts inside diffusion models [120]. Prior to diffusion models, concept-grounded interpretability has been deployed to earlier generative architectures through concept-bottleneck models, which require human

intervention at training time [121], and post-hoc detectors that retrofit concept supervision [122]. However, both of these approaches require human-defined concepts and hence inherently miss broader trends that the user does not explicitly register.

Aggregate metrics (e.g., precision, recall, density, and coverage [35, 123]) and latent density scores, which predict sample quality based on the model’s latent space [124], have emerged to evaluate generative image model capabilities. While these effectively uncover distributional gaps, they offer little insight into specific concepts that are under- or over-represented.

Together, these limitations motivate a scalable, unsupervised framework that can systematically identify and quantify concept-level failure modes in generative image models [125].

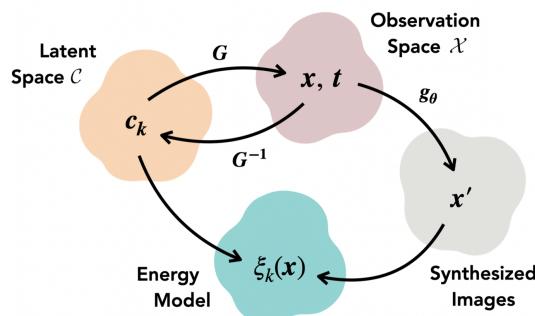
### 3 Formalizing Conceptual Blindspots in Generative Models

We begin by formalizing the notion of *conceptual blindspots*: systematic discrepancies between the conceptual content of natural images and that of model-generated outputs. This formulation enables us to derive principled, quantitative measures that characterize which concepts are under or over represented by a generative model relative to its data distribution. The process is illustrated in Fig. 2. While we rely on standard assumptions in this pursuit [126–130], empirically we find meaningful phenomenology is elicited even when these assumptions are violated.

**Definition 1 (Data-Generating Process).** Let  $\mathcal{C} \subset \mathbb{R}^K$  denote a latent space with a Boltzmann prior  $p(\mathbf{c}) = \exp(-E(\mathbf{c}))Z^{-1}$ , where  $E(\cdot)$  denotes an energy function that linearly decomposes over individual latents and  $Z$  is the corresponding partition function, i.e.,  $E(\mathbf{c}) = \sum_k E(c_k)$  and hence  $p(\mathbf{c}) = \prod_k p_k(c_k)$ . A data-generating process (DGP) is an invertible function  $G : \mathcal{C} \rightarrow \mathcal{X}$  that maps the latents  $\mathbf{c} \in \mathcal{C}$  to observations  $\mathbf{x} \in \mathcal{X}$ , i.e.,  $\mathbf{x} = G(\mathbf{c})$ .

For notational simplicity, we use  $p(\cdot)$  to denote both the latent density  $p(\mathbf{c})$  and its push-forward to image space  $p(\mathbf{x})$ , where  $\mathbf{x} = G(\mathbf{c})$ . This is justified by the invertibility of  $G$ , which induces a valid distribution over  $\mathcal{X}$  via the change of variables formula. In essence, the individual dimensions of the latent space reflect the **Concepts** underlying the data-distribution  $P_{\mathcal{X}}$ , defined over some observation space of images  $\mathcal{X}$ . For example, different latents may correspond to concepts like color, shape, size, location, and so on [131, 132]. We also let the data-generating process associate a text-description  $t \in \mathcal{T}$  with any image sampled from the data distribution, but do not explicitly model it. These text descriptions can then be used to train a text-conditioned **Generative image model**  $g_{\theta}$ , with parameters  $\theta$ , on a set of image-text pairs to map a noise signal  $\eta \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and a text-description of the scene  $t$  to produce images  $\mathbf{x}$  illustrating the latter.

To define conceptual blindspots in the model  $g_{\theta}$ , we must assess the probability mass assigned to a concept by the data-generating process, comparing it to the mass assigned by the model. To this end, we take an evaluation set of natural image-text pairs  $(\mathcal{D}_{\mathcal{X}}, \mathcal{D}_{\mathcal{T}})$  and define a set of synthetically generated images  $\mathcal{D}'_{\mathcal{X}}$  using the text descriptions. To estimate the probability of occurrence of a concept, we consider an **Energy model**  $\xi : \mathcal{X} \rightarrow \mathbb{R}^d$  that maps images to a  $d$ -dimensional vector, where each dimension associates a scalar representing the energy in the  $k^{\text{th}}$  concept, i.e.,  $\xi_k(\mathbf{x}) = E(c_k)$ . These estimates are feasible because we assume the data-generating process is invertible. Correspondingly, the probability the data-generating process associates with the image  $\mathbf{x}$  can then be defined as  $p_k(\mathbf{x}) \propto \exp(-\sum_k \xi_k(\mathbf{x}))$  (where  $Z_{\xi}$  is the partition function), hence yielding a population-level estimate  $p_k(\mathcal{D}_{\mathcal{X}}) = \prod_{\mathbf{x} \in \mathcal{D}_{\mathcal{X}}} p_k(\mathbf{x})$ . Using this and the sigmoid map  $\sigma(\cdot)$ , we define below the average energy difference in the  $k^{\text{th}}$  concept between the datasets  $\mathcal{D}_{\mathcal{X}}, \mathcal{D}'_{\mathcal{X}}$ .



**Figure 2: Intuitive mapping of our framework.** Latent concepts  $c_k \in \mathcal{C}$  are mapped to observations  $(\mathbf{x}, t)$  through the (invertible) DGP.  $S_{\theta}$  generates images  $\mathbf{x}'$  based on  $t$ . The energy model  $\xi_k(\mathbf{x})$  extracts concept representations from both  $\mathbf{x}$  and  $\mathbf{x}'$ , enabling comparison of concept distributions to identify blindspots.

**Definition 2 (Energy Difference).** Let  $\mathbf{x} \in \mathcal{D}_{\mathcal{X}}$  denote a real image sampled from the data-generating process  $\mathbf{G}$ , and let  $\mathbf{x}' \in \mathcal{D}'_{\mathcal{X}}$  be a synthetic image generated by the model  $\mathbf{g}_{\theta}$ . Let  $\xi_k : \mathcal{X} \rightarrow \mathbb{R}$  denote the energy assigned to the  $k^{\text{th}}$  concept by the energy model  $\xi$ . We define the energy difference for concept  $k$  as:

$$\begin{aligned}\delta_{\mathbf{g}_{\theta} \leftrightarrow \mathbf{G}}(k) &= \sigma(\mathbb{E}_{\mathbf{x}'}[\xi_k(\mathbf{x}')] - \mathbb{E}_{\mathbf{x}}[\xi_k(\mathbf{x})]) \\ &= \frac{p_k(\mathcal{D}'_{\mathcal{X}})}{p_k(\mathcal{D}_{\mathcal{X}}) + p_k(\mathcal{D}'_{\mathcal{X}})},\end{aligned}\quad (1)$$

where the expectations are taken over  $\mathcal{D}'_{\mathcal{X}}$  and  $\mathcal{D}_{\mathcal{X}}$ , respectively, and  $p_k(\mathcal{D}) \propto \exp(-\sum_{\mathbf{x} \in \mathcal{D}} \xi_k(\mathbf{x}))$  denotes the unnormalized conceptual probability mass of dataset  $\mathcal{D}$  under concept  $k$ .

Thus, the energy difference in the  $k^{\text{th}}$  concept describes the ratio of the probability a concept occurs in a set of observations (here,  $\mathcal{D}'_{\mathcal{X}}$ ) compared to a baseline dataset (here,  $\mathcal{D}_{\mathcal{X}}$ ). Based on this measure, we can now define conceptual blindspots as follows.

**Definition 3 (Suppressed / Exaggerated Conceptual Blindspots).** Given a generative image model  $\mathbf{g}_{\theta}$ , we say, compared to the data-generating process  $\mathbf{G}$ ,  $c_k$  is a suppressed conceptual blindspot in the model if  $\delta_{\mathbf{g}_{\theta} \leftrightarrow \mathbf{G}}(k) < \lambda_{\min}$  and exaggerated if  $\delta_{\mathbf{g}_{\theta} \leftrightarrow \mathbf{G}}(k) > \lambda_{\max}$ .

Overall, we define a conceptual blindspot as a concept whose likelihood of occurrence in generated images deviates markedly, either through suppression or exaggeration, from its prevalence under the data-generating process. Suppressed concepts exhibit disproportionately low activation (e.g.,  $\delta(k) < \lambda_{\min}$ ), whereas exaggerated concepts are overrepresented (e.g.,  $\delta(k) > \lambda_{\max}$ ). Throughout our analysis, we adopt threshold values of  $\lambda_{\min} = 0.1$  and  $\lambda_{\max} = 0.9$  to isolate these regimes.

We also note this definition is related to the idea of “mode collapse” studied in past work (e.g., see Bau et al. [39]): the difference is in the granularity at which the analysis is performed. Specifically, mode collapse focuses on exaggerated / suppressed odds of generating *entire images*, while we focus on changed odds of specific concepts. For example, if a model fails to produce images of an object with a white background, we say this concept is a suppressed conceptual blindspot.

## 4 Method: Operationalizing the Definition of Conceptual Blindspots

We next discuss our pipeline for identifying conceptual blindspots in a generative model  $\mathbf{g}_{\theta}$ . As per Sec. 3, the salient objects we need for this are (i) a set of images sampled from  $\mathbf{g}_{\theta}$  that allow comparison with the ground-truth generative process, and (ii) an energy model which enables said comparison. Below, we use  $\|\cdot\|_F$  to denote the Frobenius norm and  $\|\cdot\|_0$  to denote the number of non-zero entries (the  $\ell_0$  pseudo-norm). For a vector or matrix  $\mathbf{X}$ ,  $\mathbf{X} \geq 0$  implies element-wise non-negativity. For  $n > 0$ , we let  $[n] := \{1, \dots, n\}$ , and denote the  $i$ -th row of a matrix  $\mathbf{A}$  by  $\mathbf{A}_i$ .

**From Prompts to Latent Representations.** To identify conceptual blindspots in a model  $\mathbf{g}_{\theta}$ , we compare a dataset  $\mathcal{D}_{\mathcal{X}}$  of image-caption pairs  $(\mathbf{x}, \mathbf{t})$  sampled from the data-generating process  $\mathbf{G}$  and their synthetic counterparts sampled from the generative model  $\mathbf{g}_{\theta}$  using the text descriptions.

Specifically, given  $\mathbf{t}$ , we synthesize a counterpart image  $\mathbf{x}'$  using a pretrained text-to-image generator  $\mathbf{g}_{\theta} : \mathcal{T} \rightarrow \mathcal{X}$ , implemented as a denoising diffusion probabilistic model (DDPM) [78, 86, 87]. Sampling occurs in latent space via a reverse trajectory  $(\gamma_t)_{t=0}^T$ :

$$\gamma_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \gamma_{t-1} = \frac{1}{\sqrt{\alpha_t}} (\gamma_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \varepsilon_{\theta}(\gamma_t, t, \mathbf{c})) + \sigma_t \eta_t, \quad \text{and} \quad \eta_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

where  $\alpha_t \in (0, 1)$  and  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$  follow the standard cosine noise schedule. The final latent  $\gamma_0$  is decoded via a pretrained VAE to yield the synthetic image  $\mathbf{x}' = \text{VAE}(\gamma_0)$ . For the remainder of the paper, we treat  $\mathbf{g}_{\theta}$  as a black box that maps prompts to images:  $\mathbf{t} \mapsto \mathbf{x}'$ .

**Defining the Energy Model.** Building on prior work that shows the ability of self-supervised learning methods to invert the data-generating process and identify the energy function underlying it up to linear transformations [126, 128, 133, 134], we use DINOv2 [135]—a state-of-the-art self-supervised model—for our analysis. Under the expectation that the number of concepts underlying the DGP is larger than the dimensionality of the model’s feature space [51, 136], we train sparse autoencoders (SAEs) on its features to identify subspaces corresponding to these concepts [50, 56],

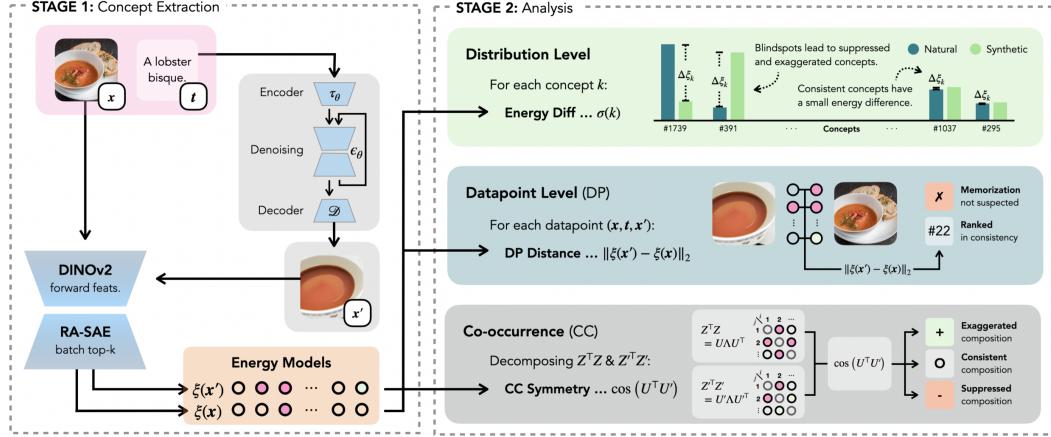


Figure 3: **Concept Extraction Pipeline.** For a triplet  $(x, t, x')$ , the concepts in  $x$  and  $t$  are extracted by obtaining each image’s DINOv2 features, which are further processed by a RA-SAE into sparse concept embeddings  $\xi(x)$  and  $\xi(x')$ , respectively. In particular,  $\xi_k(x) = E(c_k)$  holds the energy in the  $k^{\text{th}}$  concept.

137–139]. The intuition here is that if the concepts underlying the generative process are modeled via approximately orthogonal directions by DINOv2 (as assumed in our independence constraint in Def. 1), then an SAE should be able to isolate these concepts along individual dimensions in its latent space [136]. The activation associated by the SAE to a dimension will serve as our approximation of the ground-truth energy function assigned to the concept modeled by that dimension.

Formally, using  $f : \mathcal{X} \rightarrow \mathbb{R}^d$  to denote our feature extraction module (i.e., the DINOv2 model), we extract features  $a = f(x) \in \mathbb{R}^d$  from both natural and synthetic images from datasets  $\mathcal{D}_{\mathcal{X}}, \mathcal{D}'_{\mathcal{X}}$ . Assuming the count of image-text pairs is  $n$ , we stack the real and generated features into matrices  $A, A' \in \mathbb{R}^{n \times d}$ . We then decompose each feature vector into a sparse combination of learned concept atoms using an SAE. Specifically, let  $D \in \mathbb{R}^{d \times K'}$  denote a dictionary of  $K'$  concept vectors, and let  $\Psi : \mathbb{R}^d \rightarrow \mathbb{R}^{K'}$  be the SAE encoder that maps input features to sparse codes. Applying  $\Psi$  row-wise yields the matrix of activations  $Z = \Psi(A) \in \mathbb{R}^{n \times K'}$ , where each row  $z_i = \Psi(a_i)$  represents the concept decomposition of an image. The decoder reconstructs features via  $ZD^\top$ , and the SAE is trained to minimize the reconstruction error subject to sparsity and non-negativity:

$$\min_{\Psi, D} \|A - \Psi(A)D^\top\|_F^2 \quad \text{s.t.} \quad \Psi(A) \geq 0, \|\Psi(A)_i\|_0 \ll K' \quad \forall i \in [n]. \quad (2)$$

Vanilla SAEs often drift toward arbitrarily oriented dictionaries, making downstream analyses highly sensitive to the random seed. To mitigate this instability and make our study reproducible and independent of the seed, we employ the *Archetypal* SAE (RA-SAE) [56] on a TOP-K sparsity constraint [137]. RA-SAE constrains the dictionary  $D$  to be a convex combination of training data. Specifically, we write  $D = WA$  with  $W \in \Omega_{K', n}$ , the set of row-stochastic matrices in  $\mathbb{R}^{K' \times n}$ :

$$\Omega_{K', n} := \left\{ W \in \mathbb{R}^{K' \times n} \mid W \geq 0, W\mathbf{1} = \mathbf{1} \right\}. \quad (3)$$

Thus every atom  $D_i$  lies in the convex hull of the data  $\text{conv}(A)$ , and any reconstruction  $ZD^\top$  resides inside the conic hull of the data cone( $A$ ). This ensures learned concepts remain faithful to the support of the data distribution [56]. Once trained, the SAE provides a consistent set of sparse codes:  $Z$  for real images and  $Z'$  for their generated counterparts. *These codes capture the same prompt-conditioned visual semantics in terms of shared, interpretable concepts, with the activation value of the concept serving as energy values for our analysis of conceptual blindspots.* In summary then, our method defines a structured pipeline that, given a prompt and its associated real image  $(t, x)$ , produces two sparse concept vectors  $(z, z')$ , enabling direct comparison of the real and generated visual content in a common conceptual basis.

This summarizes our full pipeline: starting from a captioned image  $(t, x)$ , we synthesize a counterpart  $x'$  and map both images into a shared, sparse concept space via a vision encoder and a trained SAE, yielding interpretable representations  $(z, z')$  that will serve as the foundation for evaluating conceptual shifts induced by the generative process.

## 5 Results

We analyze four generative image models trained on LAION-5B—SD 1.5, SD 2.1, PixArt, and Kandinsky—using  $|\mathcal{X}| = 10,000$  image-text pairs and their corresponding generations (Fig. 4). Our analysis spans three levels (Fig. 3): a ● **distribution-level** evaluation reveals suppressed and exaggerated concepts; a ● **datapoint-level** analysis surfaces failures tied to ambiguity, omission, and memorization; and a ● **compositional** analysis uncovers subtle distortions in concept co-occurrence geometry.

Our core contribution is an interactive exploratory tool, shown in Fig. 5. The tool is a web-based interface built around a UMAP projection of concept representations, enabling visualization and comparison of concept-level energy differences. It is publicly available at <https://conceptual-blindspots.github.io/>, along with pre-computed energy difference data for the four models evaluated in this work (SD 1.5, SD 2.1, PixArt, and Kandinsky). All subsequent analyses in this paper are derived from insights enabled by this tool. Its primary functionalities, which support these analyses, include:

- **Contrast different models and architectures.** For each evaluated model, the tool provides a UMAP visualization spanning all 32,000 concepts from the RA-SAE. Each scatter point represents an individual concept, color-coded by its energy difference.
- **Inspect concepts.** Each concept has a card with key statistics, representative real and generated images ( $x, x'$ ), and visualized co-occurrence patterns.
- **Explore blindspots.** Beyond the UMAP and per-concept views, the tool features global rankings of suppressed and exaggerated blindspots, helping to highlight the most notable conceptual blindspots.

Rather than exhaustively studying one phenomenon, we present high-level findings that highlight the tool’s versatility and enable broader, customizable exploration.



Figure 4: **Representative image-caption pairs** from LAION-5B with matching images generated with the same prompt by SD 1.5, SD 2.1, PixArt, and Kandinsky.

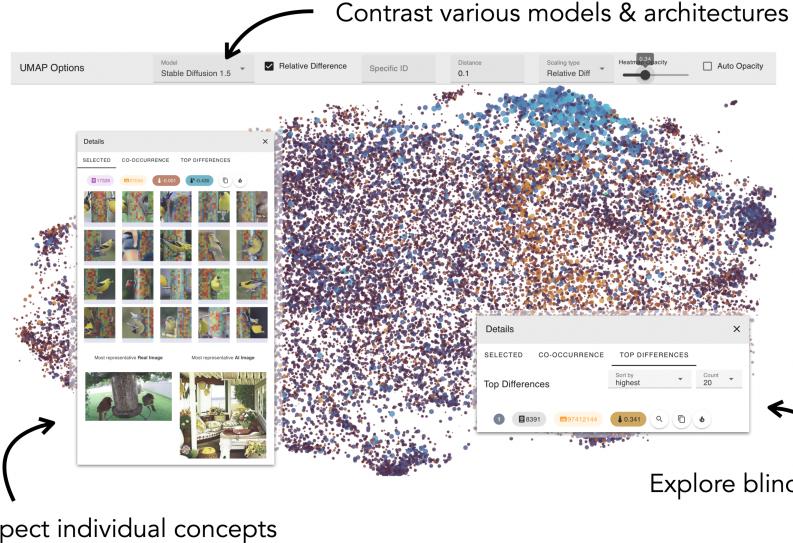


Figure 5: **Overview of the Exploratory Tool.** The web interface displays a UMAP projection for each evaluated model, where each dot represents a concept, color-coded by its energy difference. When a concept is selected, a detail panel presents illustrative images, statistics, and the most representative natural and generated images  $x$  and  $x'$ . An ordered list of the concept’s co-occurrences is shown alongside global rankings of blindspots.

## 5.1 Generative Image Models Cannot Produce Unseen Concepts

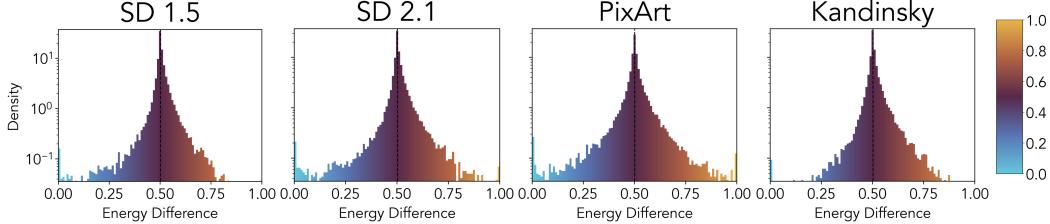


Figure 6: ● **Concept Energy Distribution.** Log-scale histograms of energy differences  $\delta(k)$  across 32,000 concepts, comparing the natural and synthesized distributions for each evaluated model. Values left of zero represent suppressed concepts (under-represented); values right of zero represent exaggerated concepts (over-represented).

To assess disparities between the generative models and the natural image distribution, we begin by analyzing the marginal energy difference  $\delta(k)$  across 32,000 concepts learned using RA-SAE. As defined in Sec. 4, this quantity reflects the relative prevalence of each concept in the synthesized versus natural image sets. A value of  $\delta(k) < 0.1$  indicates that concept  $k$  is under-represented (suppressed) in the generated images, while  $\delta(k) > 0.9$  indicates over-representation (exaggerated). Fig. 6 presents the distribution of  $\delta(k)$  for each of the four evaluated models. Across all models, we observe heavy-tailed histograms with substantial mass on both extremes, suggesting systematic discrepancies in concept coverage. Notably, the left tail—corresponding to suppressed concepts—is denser and longer than the right, indicating a consistent tendency of concept suppression. This asymmetry is reflected in the negative skewness of the distributions: Skewness =  $-0.54$  for SD 2.1,  $-0.40$  for both SD 1.5 and PixArt, and  $-0.23$  for Kandinsky.

We also note that while all models exhibit both suppressed and exaggerated concepts, their specific profiles differ. For instance, PixArt shows a wider spread, suggesting a more suppressed concept distribution. Nevertheless, the consistent left-skew in all distributions underscores a common tendency toward concept omission, though the specific characteristics of this behavior require further analysis, which we explore in the next Sections.

## 5.2 Structure and Specificity of Conceptual Blindspots

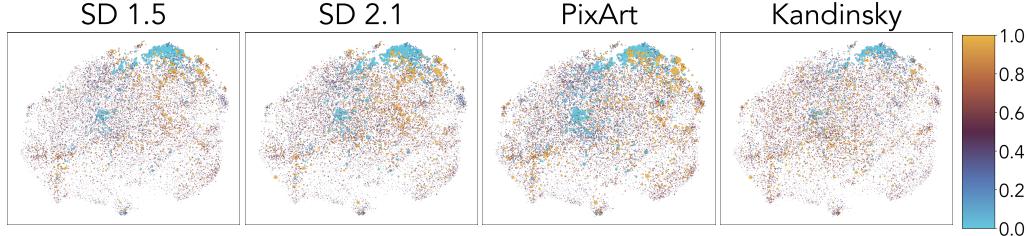


Figure 7: ● **Structure of Concept Energy Differences.** UMAP visualizations of 32,000 concepts, colored according to their energy difference  $\delta(k)$  between the natural and synthesized distributions. Clusters reveal patterns of conceptual blindspots, with suppressed concepts on the blue end and exaggerated ones on the yellow end.

While the previous section quantified marginal discrepancies in concept frequency, here we investigate their global structure by embedding the full set of 32,000 concepts into two dimensions using UMAP on the sparse codes, coloring the concepts by their  $\delta(\cdot)$  values. As shown in Fig. 7, distinct clusters of concepts emerge across all models. These clusters often correspond to contiguous regions of conceptual blindspots, especially for suppressed (blue) concepts, suggesting that blindspots are quite structured—reflecting shared biases in either training distributions or architectural priors. To assess the consistency of these blindspot patterns across models, we further analyze both the magnitude and structure of concept-level  $\delta(\cdot)$  values. Fig. 8 presents scatter plots and pairwise Pearson correlation coefficients between the  $\delta(k)$  vectors of SD 1.5 and all other models. As expected, SD 1.5 and 2.1 exhibit strong correlation ( $r = 0.82$ ), reflecting their shared architectural and training pipelines.

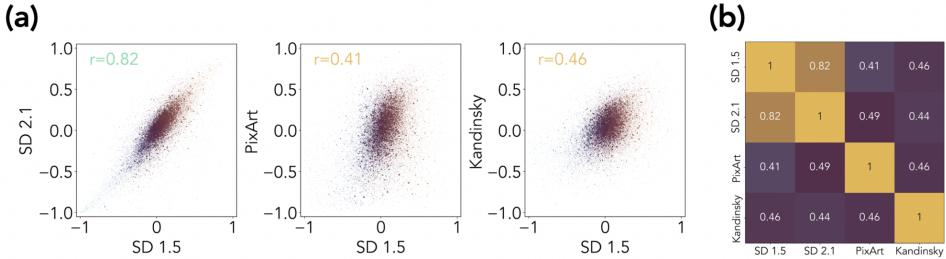


Figure 8: ● **Cross-Model Concept Energy Correlation.** Pairwise scatter plots of  $\delta(k)$  across all four evaluated models, with Pearson correlation coefficients reported top left. Strong alignment between SD 1.5 and SD 2.1 contrasts with weaker correlations among other architectures, indicating model-specific blindspots. ● **Correlation Matrix of Conceptual Blindspots.** Heatmap of pairwise Pearson correlation coefficients for  $\delta$  between all models, quantifying the degree of shared conceptual blindspots across these models.

In contrast, their correlations with PixArt and Kandinsky are substantially lower— $r = 0.41$  and  $r = 0.46$ , respectively—indicating that these models emphasize different regions of the conceptual space.

Overall, the analysis above reveals that while some blindspots are universally shared—likely due to properties of the dataset—others are highly model-specific, emerging from idiosyncrasies in training dynamics or model capacity. This motivates the need to identify and study both blindspots that are shared across models and ones that are unique to specific models in subsequent sections.

### 5.3 Qualitative Examples of Blindspots

We next visualize specific examples of both suppressed and exaggerated blindspots to gauge what concepts fall under these regimes. Specifically, in Fig. 9a we show an example of a conceptual blindspot suppressed by all models—we find all evaluated models fail to reproduce the concept `solid white on documents`. As can be seen in the figures, despite the caption explicitly referencing this concept, none of the generated images reflect the intended visual semantics, suggesting that this region of the concept space is systematically under-sampled across models. Meanwhile, Fig. 9b highlights a model-specific blindspot: the concept `pan` is accurately captured by three models, yet conspicuously missing from generations produced by Kandinsky. This reinforces the findings from Sec. 5.2, where cross-model agreement was found to be high in some cases but limited in others.

Conversely, in Fig. 10 we present a case of exaggeration, where the concept `shadow under animal` is overly emphasized in generated images. While shadows are mildly plausible, their consistent and pronounced rendering across models, relative to the more nuanced and variable occurrences in natural images, suggests an overactive prior. Interestingly, despite attempts at finding concepts that are uniquely exaggerated by a specific model, we did not find any clear examples—this suggests exaggerations are approximately universal.

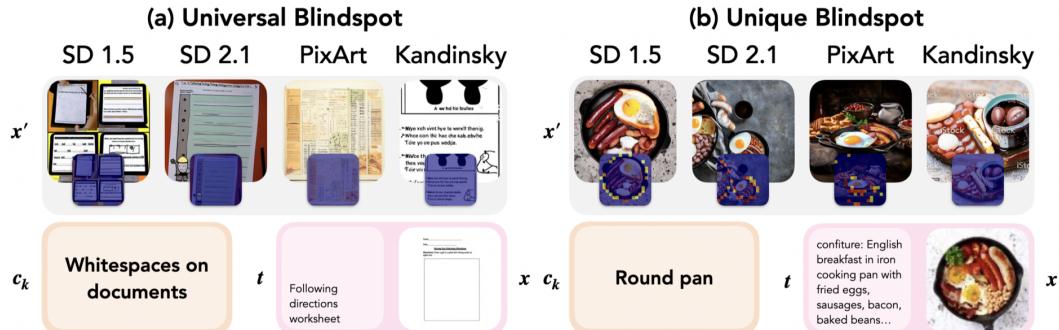


Figure 9: ● **Examples of Suppressed Conceptual Blindspots.** The natural images  $x$ , representative  $c_k$  and  $t$ , shown alongside four synthesized images  $x'$ , generated using  $S_\theta$ . The universal blindspot is present in all evaluated models; the unique blindspot is only present in Kandinsky.

Overall, the examples above concretely demonstrate how conceptual blindspots manifest in generated outputs, illustrating that our energy-based diagnostic can surface both shared and model-specific failure modes. Notably, it enables the identification of surprising model limitations—such as the consistent failure to reproduce clear or solid background elements, like whitespaces on documents, across all models. This raises the possibility that certain failure patterns may stem from architectural constraints or training data biases that transcend individual model idiosyncrasies.

While these aggregate-level analyses are informative, they invite a deeper question: do these blindspots emerge only in the aggregate across many samples, or do they manifest themselves even at the level of individual datapoints? This finer-grained perspective allows us to probe the mechanisms behind blindspots more directly—uncovering cases of prompt misinterpretation, latent memorization, or both.

#### 5.4 Investigating Datapoint-level Energy Difference: From Incongruent to Memorized Images

To move beyond population-level statistics, we examine individual natural vs. generated image pairs for which the  $\delta(\cdot)$  values averaged across all concepts exhibit the largest and smallest differences. This analysis aids easy understanding of model success and failures, latter of which we find often arises from prompt ambiguity or memorization artifacts. For example, Fig. 11a shows instances with near-zero difference in average  $\delta(\cdot)$  values. In these cases, the generated images are conceptually indistinguishable from the original. However, qualitative inspection clearly shows this happen not because the model faithfully captures the prompt semantics, but from pure replication of memorized templates: we see repetitive visual structures (e.g., outlines of clothing or object arrangements), indicating that the model may be copying from overly frequent patterns in the training data. By contrast, Fig. 11b illustrates samples that are among the largest  $\delta(\cdot)$  values. These indicate significant conceptual divergence between the synthesized and natural image. While some of these discrepancies can be attributed to underspecified or noisy captions, others reveal genuine blindspots: the prompt describes a clear concept faithfully present in  $\mathcal{D}_{\mathcal{X}}$ , yet the model fails to realize it in  $\mathcal{D}'_{\mathcal{X}}$ . This failure suggests that even when language grounding is adequate, certain concepts fall outside the model’s generative abilities.



Figure 10: ● **Example of an Exaggerated Conceptual Blindspot.** Four synthesized images  $\mathcal{D}'_{\mathcal{X}}$ , archetypal for the concept shadow under animal, alongside the corresponding natural image  $\mathcal{D}_{\mathcal{X}}$  and caption  $t$  used to generate it. Note the accentuated, nearly solid black shadows underneath the animals. We deem this blindspot universal since it is present in all evaluated models.



Figure 11: ● **Datapoint-level Conceptual Alignment.** (a) Examples with minimal energy differences where models appear to memorize training patterns. (b) Examples with large differences where significant concept divergences due to prompt ambiguity or model limitations occur.

## 5.5 Analyzing Post-Training Effects

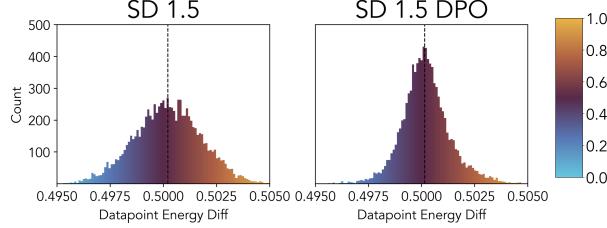


Figure 12: ● **Effect of DPO on Concept Fidelity.** Histogram of datapoint-wise energy differences between the synthesized and natural distribution of SD 1.5 models with and without DPO.

Post-training protocols, e.g., safety fine-tuning, have been argued to reduce the diversity of model generations [140]. Given our pipeline’s ability to isolate interesting differences in a model’s generations and the ground-truth DGP, we next use it to understand the effects of DPO—a popular safety fine-tuning protocol [141]. Specifically, we compare two variants of the SD 1.5 model: one trained with DPO, and one without. For each image pair  $(\mathcal{D}_X, \mathcal{D}'_X)$ , we compute the  $\ell_2$  norm of the difference between their internal concept energy vectors,  $\|\xi(\mathcal{D}'_X) - \xi(\mathcal{D}_X)\|_2$ . Fig. 12 presents a histogram of these datapoint-wise energy differences. The DPO-enhanced model exhibits both a lower median and a narrower spread, indicating more consistent distribution of generated concepts with the ground-truth DGP. This suggests that DPO may serve to regularize the model’s concept distribution, encouraging outputs that more closely reflect the semantic content of the seen inputs. While our analysis does not disentangle the specific inductive biases introduced by DPO, these results provide empirical evidence that its optimization objective, which favors human-preferred generations, indirectly promotes better match with the training distribution. In particular, it reduces *both over and under* activation of individual concepts relative to the natural baseline. These findings highlight the utility of our pipeline in characterizing the downstream effects of post-training interventions: not merely in terms of output quality, but in how they reshape the conceptual geometry of the model’s output space.

## 5.6 Conceptual Misalignment as a Function of Empirical Frequency

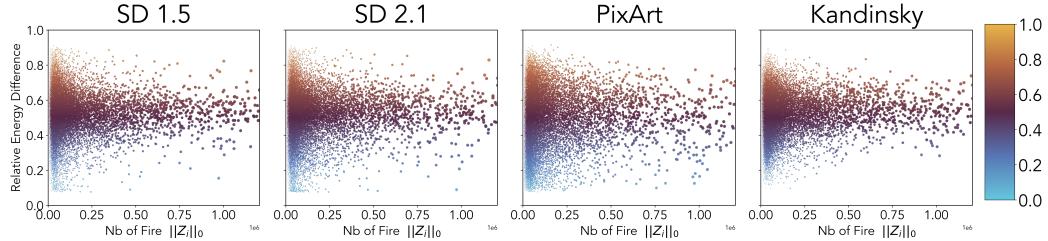


Figure 13: ● **Concept Fidelity Across Frequency Spectrum.** Scatter plots showing the relationship between concept frequency (x-axis) and the energy difference (y-axis) across four evaluated models. Each point represents a concept with size is proportional to its activation frequency.

We previously hypothesized that conceptual blindspots are not merely be architectural artifacts, but may also emerge as a direct consequence of distributional peculiarities of certain features. In this section, we empirically test this hypothesis by examining whether concepts that are rarely activated in natural images, i.e., those in the long tail of the data distribution, lead to blindspots in generative models. Specifically, we process the natural dataset  $\mathcal{D}_X$  through the trained SAE and compute, for each concept  $k$ , its empirical frequency  $\|\mathbf{Z}_{:,k}\|_0$ , where  $\mathbf{Z}_{:,i}$  is the activations of concept  $i$  across all our images. We then correlate this with the absolute energy difference observed across generated outputs. Fig. 13 visualizes this relationship for all evaluated models. We find that concepts with higher frequency in natural data tend to show lower energy discrepancies, while rare concepts—especially suppressed ones ( $\delta(k) < 0.5$ )—exhibit significant alignment errors. This suggests that many blindspots stem not from randomness or model quirks, but from systematic effects tied to long-tail concept distributions. Addressing these issues may require not just architectural changes but also strategies like data reweighting or augmentation.

## 6 Discussion

Our analysis reveals multiple conceptual blindspots in four popular generative image models. The results presented here, however, only scratch the surface: each individual finding could warrant its own dedicated investigation. Rather than delving deeply into any one of these questions, we instead showcase the versatility of our method and exploratory tool. Out of the box, they allow for a systematic identification of concepts that models struggle to generate, detection of memorization artifacts, discovery of datapoints with insufficient captions, quantification of post-training effects, and characterization of conceptual shifts across model architectures. We thus open space for follow-up work to extend the depth of analysis, scope of evaluated architectures, and inquiry into root causes of conceptual blindspots. Future work could also explore hierarchical representations of concepts in RA-SAE to allow for a more nuanced analysis.

**Limitations.** We wish to highlight several limitations of our work. By relying on DINOv2 and RA-SAE for concept extraction and representation, our approach is inherently constrained to the kinds of concepts these models capture; concepts poorly represented by them will escape our analysis. Additionally, while our sample size of 10,000 images is substantial, it may not fully capture the long tail of rare concepts, concept co-occurrence, or other compositional statistics.

## References

- [1] OpenAI. Sora: Creating video from text, 2024. URL <https://openai.com/sora/>.
- [2] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023.
- [3] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pages 8821–8831. PMLR, 2021.
- [4] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022.
- [5] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. GLIDE: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv:2112.10741*, 2021.
- [6] Chen Wang, Jiatao Gu, Xiaoxiao Long, Yuan Liu, and Lingjie Liu. Geco: Generative image-to-3d within a second. *arXiv preprint arXiv:2405.20327*, 2024.
- [7] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022.
- [8] Elad Richardson, Gal Metzer, Yuval Alaluf, Raja Giryes, and Daniel Cohen-Or. Texture: Text-guided texturing of 3d shapes. In *ACM SIGGRAPH 2023 conference proceedings*, pages 1–11, 2023.
- [9] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [10] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- [11] Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, et al. Inner monologue: Embodied reasoning through planning with language models. *arXiv preprint arXiv:2207.05608*, 2022.
- [12] Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International Conference on Machine Learning*, pages 9118–9147. PMLR, 2022.
- [13] Yun-Chun Chen, Selena Ling, Zhiqin Chen, Vladimir G Kim, Matheus Gadelha, and Alec Jacobson. Text-guided controllable mesh refinement for interactive 3d modeling. In *SIGGRAPH Asia 2024 Conference Papers*, pages 1–11, 2024.
- [14] Lei Zhong, Yiming Xie, Varun Jampani, Deqing Sun, and Huaizu Jiang. Smoodi: Stylized motion diffusion model. In *European Conference on Computer Vision*, pages 405–421. Springer, 2024.
- [15] Yawar Siddiqui, Antonio Alliegro, Alexey Artemov, Tatiana Tommasi, Daniele Sirigatti, Vladislav Rosov, Angela Dai, and Matthias Nießner. Meshgpt: Generating triangle meshes with decoder-only transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19615–19625, 2024.
- [16] Gary Marcus, Ernest Davis, and Scott Aaronson. A very preliminary analysis of dall-e 2. *arXiv preprint arXiv:2204.13807*, 2022.
- [17] Ángel Alexander Cabrera, Abraham J Druck, Jason I Hong, and Adam Perer. Discovering and validating ai errors with crowdsourced failure reports. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW2):1–22, 2021.

- [18] Rebecca Heigl. Generative artificial intelligence in creative contexts: a systematic review and future research agenda. *Management Review Quarterly*, pages 1–38, 2025.
- [19] Qihao Liu, Adam Kortylewski, Yutong Bai, Song Bai, and Alan Yuille. Discovering failure modes of text-guided diffusion models via adversarial search. *arXiv preprint arXiv:2306.00974*, 2023.
- [20] Colin Conwell, Rupert Tawiah-Quashie, and Tomer Ullman. Relations, negations, and numbers: Looking for logic in generative text-to-image models. *arXiv preprint arXiv:2411.17066*, 2024.
- [21] Wenquan Lu, Yufei Xu, Jing Zhang, Chaoyue Wang, and Dacheng Tao. HandRefiner: Refining malformed hands in generated images by diffusion-based conditional inpainting. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 7085–7093, 2024.
- [22] Supreeth Narasimhaswamy, Uttaran Bhattacharya, Xiang Chen, Ishita Dasgupta, Saayan Mitra, and Minh Hoai. HanDiffuser: Text-to-image generation with realistic hand appearances. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2468–2479, 2024.
- [23] Qilong Zhangli, Jindong Jiang, Di Liu, Licheng Yu, Xiaoliang Dai, Ankit Ramchandani, Guan Pang, Dimitris N Metaxas, and Praveen Krishnan. Layout-agnostic scene text image synthesis with diffusion models. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7496–7506. IEEE Computer Society, 2024.
- [24] Forouzan Fallah, Maitreya Patel, Agneet Chatterjee, Vlad I Morariu, Chitta Baral, and Yezhou Yang. TextInVision: Text and prompt complexity driven visual text generation benchmark. *arXiv:2503.13730*, 2025.
- [25] Yuefan Cao, Xuyang Guo, Jiayan Huo, Yingyu Liang, Zhenmei Shi, Zhao Song, Jiahao Zhang, and Zhen Zhuang. Text-to-image diffusion models cannot count, and prompt refinement cannot help. *arXiv preprint arXiv:2503.06884*, 2025.
- [26] Colin Conwell and Tomer Ullman. Testing relational understanding in text-guided image generation. *arXiv preprint arXiv:2208.00005*, 2022.
- [27] Mahzarin R Banaji and Anthony G Greenwald. *Blindspot: Hidden biases of good people*. Bantam, 2016.
- [28] George Stein, Jesse Cresswell, Rasa Hosseinzadeh, Yi Sui, Brendan Ross, Valentin Villecroze, Zhaoyan Liu, Anthony L Caterini, Eric Taylor, and Gabriel Loaiza-Ganem. Exposing flaws of generative model evaluation metrics and their unfair treatment of diffusion models. *Advances in Neural Information Processing Systems*, 36:3732–3784, 2023.
- [29] Boheng Wang, Yunhuai Zhu, Liuqing Chen, Jingcheng Liu, Lingyun Sun, and Peter Childs. A study of the evaluation metrics for generative images containing combinational creativity. *AI EDAM*, 37:e11, 2023.
- [30] 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. *Advances in neural information processing systems*, 30, 2017.
- [31] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7514–7528, 2021.
- [32] Mischa Dombrowski, Weitong Zhang, Sarah Cechnicka, Hadrien Reynaud, and Bernhard Kainz. Image generation diversity issues and how to tame them. *arXiv:2411.16171*, 2024.
- [33] Jaehui Hwang, Junghyuk Lee, and Jong-Seok Lee. Anomaly score: Evaluating generative models and individual generated images based on complexity and vulnerability. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8754–8763, 2024.

- [34] Lucas Theis, Aäron van den Oord, and Matthias Bethge. A note on the evaluation of generative models. arXiv:1511.01844, 2015.
- [35] Muhammad Ferjad Naeem, Seong Joon Oh, Youngjung Uh, Yunjey Choi, and Jaejun Yoo. Reliable fidelity and diversity metrics for generative models. In *International conference on machine learning*, pages 7176–7185. PMLR, 2020.
- [36] Vitali Petsiuk, Alexander E Siemann, Saisamrit Surbehera, Zad Chin, Keith Tyser, Gregory Hunter, Arvind Raghavan, Yann Hicke, Bryan A Plummer, Ori Kerret, et al. Human evaluation of text-to-image models on a multi-task benchmark. arXiv:2211.12112, 2022.
- [37] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36:15903–15935, 2023.
- [38] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. 2023.
- [39] David Bau, Jun-Yan Zhu, Jonas Wulff, William Peebles, Hendrik Strobelt, Bolei Zhou, and Antonio Torralba. Seeing what a GAN cannot generate. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4502–4511, 2019.
- [40] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2013.
- [41] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. *Proceedings of the International Conference on Machine Learning (ICML)*, 2017.
- [42] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [43] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. *Advances in Neural Information Processing Systems (NIPS)*, 2018.
- [44] Amirata Ghorbani, Abubakar Abid, and James Zou. Interpretation of neural networks is fragile. *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2017.
- [45] Peter Hase and Mohit Bansal. Evaluating explainable ai: Which algorithmic explanations help users predict model behavior? *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*, 2020.
- [46] 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). *Proceedings of the International Conference on Machine Learning (ICML)*, 2018.
- [47] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [48] Thomas Fel, Agustin Picard, Louis Bethune, Thibaut Boissin, David Vigouroux, Julien Colin, Rémi Cadène, and Thomas Serre. Craft: Concept recursive activation factorization for explainability. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [49] Matthew Kowal, Achal Dave, Rares Ambrus, Adrien Gaidon, Konstantinos G Derpanis, and Pavel Tokmakov. Understanding video transformers via universal concept discovery. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.

- [50] Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. *ArXiv e-print*, 2023.
- [51] Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023.
- [52] Thomas Fel, Victor Boutin, Mazda Moayeri, Remi Cadene, Louis Bethune, Mathieu Chalvidal, and Thomas Serre. A holistic approach to unifying automatic concept extraction and concept importance estimation. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- [53] David Chanin, James Wilken-Smith, Tomas Dulka, Hardik Bhatnagar, and Joseph Bloom. A is for absorption: Studying feature splitting and absorption in sparse autoencoders. *ArXiv e-print*, 2024.
- [54] Martin Wattenberg and Fernanda B Viegas. Relational composition in neural networks: A survey and call to action. *ArXiv e-print*, 2024.
- [55] Usha Bhalla, Suraj Srinivas, Asma Ghandeharioun, and Himabindu Lakkaraju. Towards unifying interpretability and control: Evaluation via intervention. *ArXiv e-print*, 2024.
- [56] Thomas Fel, Ekdeep Singh Lubana, Jacob S Prince, Matthew Kowal, Victor Boutin, Isabel Papadimitriou, Binxu Wang, Martin Wattenberg, Demba Ba, and Talia Konkle. Archetypal SAE: Adaptive and stable dictionary learning for concept extraction in large vision models. 2025.
- [57] Arpita Chowdhury, Dipanjyoti Paul, Zheda Mai, Jianyang Gu, Ziheng Zhang, Kazi Sajeed Mehrab, Elizabeth G Campolongo, Daniel Rubenstein, Charles V Stewart, Anuj Karpatne, et al. Prompt-CAM: A simpler interpretable transformer for fine-grained analysis. arXiv:2501.09333, 2025.
- [58] David Bau, Jun-Yan Zhu, Hendrik Strobelt, Bolei Zhou, Joshua B Tenenbaum, William T Freeman, and Antonio Torralba. GAN dissection: Visualizing and understanding generative adversarial networks. 2018.
- [59] Raphael Tang, Lingqiang Liu, Akshat Pandey, Zhiying Jiang, Gefei Yang, Karun Kumar, Pontus Stenetorp, Jimmy Lin, and Ferhan Ture. What the DAAM: Interpreting stable diffusion using cross attention. 2022.
- [60] Ali Borji. Qualitative failures of image generation models and their application in detecting deepfakes. *Image and Vision Computing*, 137:104771, 2023.
- [61] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with CLIP latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- [62] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- [63] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. 2023.
- [64] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International conference on machine learning*, pages 8162–8171. PMLR, 2021.
- [65] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022.
- [66] Haoyu Lu, Guoxing Yang, Nanyi Fei, Yuqi Huo, Zhiwu Lu, Ping Luo, and Mingyu Ding. VDT: General-purpose video diffusion transformers via mask modeling. *arXiv preprint arXiv:2305.13311*, 2023.

- [67] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. ModelScope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023.
- [68] Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. AnimateDiff: Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023.
- [69] Shanchuan Lin and Xiao Yang. AnimateDiff-Lightning: Cross-model diffusion distillation. *arXiv preprint arXiv:2403.12706*, 2024.
- [70] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. CogVideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022.
- [71] Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, et al. VideoCrafter1: Open diffusion models for high-quality video generation. *arXiv preprint arXiv:2310.19512*, 2023.
- [72] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7623–7633, 2023.
- [73] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3D: High-resolution text-to-3d content creation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 300–309, 2023.
- [74] Heewoo Jun and Alex Nichol. Shap-E: Generating conditional 3d implicit functions. *arXiv preprint arXiv:2305.02463*, 2023.
- [75] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12619–12629, 2023.
- [76] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [77] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [78] Stability AI. Stable Diffusion 2.0 Release, 2022. URL <https://stability.ai/news/stable-diffusion-v2-release>.
- [79] Stability AI. Introducing Stable Diffusion 3, 2024. URL <https://stability.ai/news/stable-diffusion-3>.
- [80] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- [81] Kihyuk Sohn, Nataniel Ruiz, Kimin Lee, Daniel Castro Chin, Irina Blok, Huiwen Chang, Jarred Barber, Lu Jiang, Glenn Entis, Yuanzhen Li, et al. StyleDrop: Text-to-image generation in any style. *arXiv preprint arXiv:2306.00983*, 2023.
- [82] Zhihong Pan, Xin Zhou, and Hao Tian. Arbitrary style guidance for enhanced diffusion-based text-to-image generation. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 4461–4471, 2023.
- [83] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3836–3847, 2023.

- [84] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. DreamBooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22500–22510, 2023.
- [85] Darian Tomavsevic, Fadi Boutros, Chenhao Lin, Naser Damer, Vitomir Struc, and Peter Peer. Id-booth: Identity-consistent face generation with diffusion models. *arXiv preprint arXiv:2504.07392*, 2025.
- [86] Anton Razzhigaev, Arseniy Shakhmatov, Anastasia Maltseva, Vladimir Arkhipkin, Igor Pavlov, Ilya Ryabov, Angelina Kuts, Alexander Panchenko, Andrey Kuznetsov, and Denis Dimitrov. Kandinsky: An improved text-to-image synthesis with image prior and latent diffusion. *arXiv preprint arXiv:2310.03502*, 2023.
- [87] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart- $\alpha$ : Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023.
- [88] Stability AI. DeepFloyd IF: A Powerful Open-Source Text-to-Image Model, 2023. URL <https://stability.ai/news/deepfloyd-if-text-to-image-model>.
- [89] Black Forest Labs. Announcing Black Forest Labs, 2024. URL <https://blackforestlabs.ai/announcing-black-forest-labs/>.
- [90] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- [91] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2(3):5, 2022.
- [92] Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad, Md Mostofa Ali Patwary, Yang Yang, and Yanqi Zhou. Deep learning scaling is predictable, empirically. *arXiv preprint arXiv:1712.00409*, 2017.
- [93] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. LAION-5B: An open large-scale dataset for training next generation image-text models. *Advances in neural information processing systems*, 35:25278–25294, 2022.
- [94] Common Crawl. Common crawl corpus. <https://commoncrawl.org>, 2017.
- [95] Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Saehoon Kim. COYO-700M: Image-text pair dataset, 2022.
- [96] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual Captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of ACL*, 2018.
- [97] Abeba Birhane, Sanghyun Han, Vishnu Boddeti, Sasha Luccioni, et al. Into the LAION’s den: Investigating hate in multimodal datasets. *Advances in neural information processing systems*, 36:21268–21284, 2023.
- [98] Abeba Birhane, Sepehr Dehdashtian, Vinay Prabhu, and Vishnu Boddeti. The dark side of dataset scaling: Evaluating racial classification in multimodal models. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, pages 1229–1244, 2024.
- [99] Preethi Seshadri, Sameer Singh, and Yanai Elazar. The bias amplification paradox in text-to-image generation. *arXiv preprint arXiv:2308.00755*, 2023.

- [100] Abeba Birhane, Vinay Uday Prabhu, and Emmanuel Kahembwe. Multimodal datasets: Misogyny, pornography, and malignant stereotypes. *arXiv preprint arXiv:2110.01963*, 2021.
- [101] David Thiel. Identifying and eliminating CSAM in generative ml training data and models. *Stanford Internet Observatory, Cyber Policy Center, December*, 23:3, 2023.
- [102] Ali Shirali and Moritz Hardt. What makes ImageNet look unlike LAION. *arXiv preprint arXiv:2306.15769*, 2023.
- [103] Thao Nguyen, Samir Yitzhak Gadre, Gabriel Ilharco, Sewoong Oh, and Ludwig Schmidt. Improving multimodal datasets with image captioning. *Advances in Neural Information Processing Systems*, 36:22047–22069, 2023.
- [104] Bruno A Olshausen and David J Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 1996.
- [105] Michael Elad. Sparse and redundant representations: from theory to applications in signal and image processing. 2010.
- [106] Julien Mairal, Francis Bach, and Jean Ponce. Sparse modeling for image and vision processing. *Foundations and Trends in Computer Graphics and Vision*, 2014.
- [107] David L Donoho. Compressed sensing. *IEEE Transactions on Information Theory*, 2006.
- [108] Emmanuel J Candès, Justin Romberg, and Terence Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on Information Theory*, 2006.
- [109] Stuart Lloyd. Least squares quantization in pcm. *IEEE Transactions on Information Theory*, 1982.
- [110] Daniel D Lee and H Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 1999.
- [111] Nicolas Gillis. Nonnegative matrix factorization. 2020.
- [112] Hui Zou, Trevor Hastie, and Robert Tibshirani. Sparse principal component analysis. *Journal of Computational and Graphical Statistics*, 2006.
- [113] Michal Aharon, Michael Elad, and Alfred Bruckstein. K-svd: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on Signal Processing*, 2006.
- [114] Julien Mairal, Francis Bach, Jean Ponce, and Guillermo Sapiro. Online dictionary learning for sparse coding. *Proceedings of the International Conference on Machine Learning (ICML)*, 2009.
- [115] Rodolphe Jenatton, Guillaume Obozinski, and Francis Bach. Structured sparse principal component analysis. *International Conference on Artificial Intelligence and Statistics*, 2010.
- [116] Daniel A Spielman, Huan Wang, and John Wright. Exact recovery of sparsely-used dictionaries. *The Journal of Machine Learning Research (JMLR)*, 2012.
- [117] Jean Barbier and Nicolas Macris. Statistical limits of dictionary learning: random matrix theory and the spectral replica method. *Physical Review E*, 2022.
- [118] Vardan Papyan, Yaniv Romano, and Michael Elad. Convolutional dictionary learning via local processing. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [119] Alex Tamkin, Mohammad Taufeeque, and Noah D Goodman. Codebook features: Sparse and discrete interpretability for neural networks. *ArXiv e-print*, 2023.
- [120] Berk Tinaz, Zalan Fabian, and Mahdi Soltanolkotabi. Emergence and evolution of interpretable concepts in diffusion models. *arXiv:2504.15473*, 2025.

- [121] Akshay Kulkarni, Ge Yan, Chung-En Sun, Tuomas Oikarinen, and Tsui-Wei Weng. Interpretable generative models through post-hoc concept bottlenecks. *arXiv preprint arXiv:2503.19377*, 2025.
- [122] Mert Yuksekgonul, Maggie Wang, and James Zou. Post-hoc concept bottleneck models. *arXiv preprint arXiv:2205.15480*, 2022.
- [123] Tuomas Kynkänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved precision and recall metric for assessing generative models. *Advances in neural information processing systems*, 32, 2019.
- [124] Jingyi Xu, Hieu Le, and Dimitris Samaras. Assessing sample quality via the latent space of generative models. In *European Conference on Computer Vision*, pages 449–464. Springer, 2024.
- [125] Iro Laina, Yuki M Asano, and Andrea Vedaldi. Measuring the interpretability of unsupervised representations via quantized reverse probing. *arXiv:2209.03268*, 2022.
- [126] Julius Von Kugelgen, Yash Sharma, Luigi Gresele, Wieland Brendel, Bernhard Scholkopf, Michel Besserve, and Francesco Locatello. Self-supervised learning with data augmentations provably isolates content from style. *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [127] Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard Scholkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. *Proceedings of the International Conference on Machine Learning (ICML)*, 2019.
- [128] Roland S Zimmermann, Yash Sharma, Steffen Schneider, Matthias Bethge, and Wieland Brendel. Contrastive learning inverts the data generating process. *Proceedings of the International Conference on Machine Learning (ICML)*, 2021.
- [129] Luigi Gresele, Paul K Rubenstein, Arash Mehrjou, Francesco Locatello, and Bernhard Scholkopf. The incomplete rosetta stone problem: Identifiability results for multi-view nonlinear ica. *Uncertainty in Artificial Intelligence*, 2020.
- [130] Luigi Gresele, Julius Von Kugelgen, Vincent Stimper, Bernhard Scholkopf, and Michel Besserve. Independent mechanism analysis, a new concept? *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [131] Maya Okawa, Ekdeep S Lubana, Robert Dick, and Hidenori Tanaka. Compositional abilities emerge multiplicatively: Exploring diffusion models on a synthetic task. *Advances in Neural Information Processing Systems*, 36:50173–50195, 2023.
- [132] Core Francisco Park, Maya Okawa, Andrew Lee, Ekdeep S Lubana, and Hidenori Tanaka. Emergence of hidden capabilities: Exploring learning dynamics in concept space. *Advances in Neural Information Processing Systems*, 37:84698–84729, 2024.
- [133] Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational autoencoders and nonlinear ica: A unifying framework. *Proceedings of the International Conference on Machine Learning (ICML)*, 2020.
- [134] Aapo Hyvarinen, 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.
- [135] Maxime Oquab, Timothée Darct, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *ArXiv e-print*, 2023.
- [136] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition. *Transformer Circuits Thread*, 2022.

- [137] Leo Gao, Tom Dupre la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025.
- [138] Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Transformer Circuits Thread*, 2024.
- [139] Senthooran Rajamanoharan, Tom Lieberum, Nicolas Sonnerat, Arthur Conmy, Vikrant Varma, Janos Kramar, and Neel Nanda. Jumping ahead: Improving reconstruction fidelity with jumprelu sparse autoencoders. *ArXiv e-print*, 2024.
- [140] Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward Grefenstette, and Roberta Raileanu. Understanding the effects of rlhf on llm generalisation and diversity. *arXiv preprint arXiv:2310.06452*, 2023.
- [141] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023.
- [142] Colin McDiarmid et al. On the method of bounded differences. *Surveys in combinatorics*, 141(1):148–188, 1989.

# Appendix

## Table of Contents

---

<b>A Experimental Setup</b>	<b>23</b>
A.1 Observation Space . . . . .	23
A.2 Synthesized Images . . . . .	23
A.3 ● Distribution Level Analysis . . . . .	23
A.4 ● Datapoint Level Analysis . . . . .	24
A.5 ● Co-occurrence Analysis . . . . .	24
<b>B Custom RA-SAE Training Details</b>	<b>25</b>
<b>C Computational Resources</b>	<b>26</b>
<b>D Autointerpretability</b>	<b>27</b>
<b>E Additional Results: Model-Specific Blindspots</b>	<b>29</b>
<b>F Additional Results: Higher-order Blindspots with Compositional Discrepancy</b>	<b>31</b>
<b>G Additional Results: Blindspot Stress Testing</b>	<b>33</b>
G.1 Bird Feeder Blindspot in Kandinsky . . . . .	33
G.2 Glossy DVD Disc Blindspot in SD 1.5 . . . . .	34
<b>H Additional Examples: Synthesized Images</b>	<b>35</b>
<b>I Concentration bounds for <math>\delta</math></b>	<b>36</b>
<b>J Monotonicity and Calibration-Free Interpretation of <math>\delta</math></b>	<b>37</b>
<b>K Stability of FID Under SAE Embeddings</b>	<b>38</b>

---

## A Experimental Setup

This section details the experimental setup for our analysis of four popular generative image models: SD 1.5/2.1, Kandinsky, and PixArt, all trained on LAION-5B or its subsets/derivatives. The code is available at <https://conceptual-blindspots.github.io/>.

### A.1 Observation Space

The observation space is constructed by sampling 10,000 image-text pairs from the LAION-5B dataset [93], which serves as our domain of natural images. Due to concerns with CSAM and other unsafe content in the dataset, the original data release is no longer available. A substitute release of a subset of this dataset with additional filtering of the unsafe content, available at <https://huggingface.co/datasets/laion/relaion2B-en-research-safe>, is used.

The sampling procedure consists of: (1) loading the full LAION dataset using the Hugging Face datasets library, (2) performing validation to ensure proper URL structure and resource availability via HTTP HEAD requests, and (3) employing random sampling with replacement until reaching the target count of 10,000 valid samples. This approach, yielding  $D_G$  with  $(\mathbf{x}, t)$  tuples, ensures our observation space contains accessible image-text pairs for comparative analysis of a dataset of image URLs whose large portion has been made unavailable since original release. Additional examples of synthesized images are shown in Appendix H.

### A.2 Synthesized Images

For each of the four evaluated models, we generate a synthetic dataset  $D_{g_\theta}$  to have a one-to-one correspondence with  $D_G$ , yielding triplets  $(\mathbf{x}, \mathbf{x}', t)$ . Specifically, given the 10,000 image-text pairs  $(\mathbf{x}, t)$  from  $D_G$ , we use  $t$  to synthesize counterpart images  $\mathbf{x}'$  using each generative model  $g_\theta$ .

The synthesis process follows the standard text-to-image generation pipeline for each model architecture, implemented using the Hugging Face diffusers library, where the models are loaded at mixed precision (fp16). All synthetic images are generated at  $512 \times 512$  pixel resolution with default parameters.

**Stable Diffusion 1.5.** The checkpoint from <https://huggingface.co/benjamin-paine/stable-diffusion-v1-5> (which is a mirror of the deprecated <https://huggingface.co/ruwnayml/stable-diffusion-v1-5>) is used. Inference is performed using 50 inference steps, with the guidance scale fixed at 7.5.

**Stable Diffusion 1.5 + DPO.** The DPO variant of SD 1.5 (used in the analysis in Sec. 5.5) follows the baseline SD 1.5 implementation, but replaces the UNet component with a DPO-trained version from <https://huggingface.co/mhdang/dpo-sd1.5-text2image-v1>.

**Stable Diffusion 2.1.** The checkpoint from <https://huggingface.co/stabilityai/stable-diffusion-2-1> is used. Inference is performed using 50 inference steps, with the guidance scale fixed at 7.5.

**Kandinsky.** The checkpoint from <https://huggingface.co/kandinsky-community/kandinsky-2-1> is used. Inference is performed using 100 inference steps, with the guidance scale fixed at 4.0.

**PixArt.** The checkpoint from <https://huggingface.co/PixArt-alpha/PixArt-XL-2-1024-MS> is used. Inference is performed using 50 inference steps, with the guidance scale fixed at 7.5.

### A.3 ● Distribution Level Analysis

**Section 5.1.** We compute energy differences  $\delta(\cdot)$  across all 32,000 concepts for each evaluated model. The sigmoid transformation with temperature  $T = 0.8$  is applied during normalization. The resulting values are visualized as log-scale density histograms with 100 bins spanning  $[0, 1]$ .

**Section 5.2.** We embed the complete set of 32,000 concepts into two-dimensional space using UMAP applied to the sparse concept codes. Each point in this UMAP represents an individual concept, colored according to its energy difference  $\delta(\cdot)$ , emphasizing both suppressed and exaggerated blindspots. To quantify cross-model consistency, we compute pairwise Pearson correlation coefficients between  $\delta(\cdot)$  vectors of all model pairs, producing both scatter plots and correlation matrices. This analysis reveals whether blindspots cluster in conceptual space and identifies model-specific versus universal patterns of conceptual blindspots.

**Section 5.3.** We rank all 32,000 concepts by their energy difference  $\delta(\cdot)$ , and manually examine the extrema (both suppressed and exaggerated blindspots). For suppressed blindspots, we select concepts with  $\delta(\cdot) < 0.1$ ; for exaggerated blindspots, we choose those with  $\delta(\cdot) < 0.9$ . Presented examples are manually annotated with textual descriptions of the respective concepts through inspection of their most activating images and spatial attention patterns. We outline ongoing efforts to automate this concept interpretation in Appendix D.

**Section 5.5.** We compare 1.5 with and without DPO in the following fashion: for each image pair  $(x, x')$ , we compute the L2 norm of the difference between their concept energy vectors  $\|\xi(x') - \xi(x)\|_2$ . We apply a sigmoid transformation with temperature  $T = 0.8$  to the element-wise differences before taking their mean. This yields datapoint-wise energy differences that quantify how much each generated image deviates from its natural counterpart in concept space. Finally, these differences are visualized as overlapping histograms, contrasting both model variants.

**Section 5.6.** For each concept  $c_k$ , its empirical frequency  $\|Z_{:,k}\|_0$  (the count of non-zero activations across the natural dataset) is counted. A sigmoid normalization with temperature  $T = 0.4$  is then applied to the energy differences  $\delta(\cdot)$ . The analysis is visualized using scatter plots where the x-axis is the empirical concept frequency and the y-axis is the sigmoid-transformed energy difference. The point sizes are proportional to activation frequency and point colors are proportional to the magnitude of energy differences.

#### A.4 • Datapoint Level Analysis

**Section 5.4.** For each image pair  $(x, x')$ , we compute the L2 norm of the difference between their concept energy vectors  $\|\xi(x') - \xi(x)\|_2$ . This yields a scalar measure of conceptual divergence for each image pair. The samples are ranked by their energy differences. Minimal divergence indicate potential memorization artifacts and maximal divergence point to significant conceptual failures. This analysis enables qualitative inspection of specific failure modes.

#### A.5 • Co-occurrence Analysis

**Appendix F** For both the natural and synthesized data  $D_G$  and  $D_{g\theta}$ , concept co-occurrence patterns are analyzed through the co-activation matrix  $Z^T Z$ , which holds pair-wise correlations in concept usage. Spectral analysis is performed to examine the dominant conceptual directions using eigendecomposition. The alignment between natural and synthetic co-occurrence structures is assessed using cosine similarity heatmaps between the top-100 eigenvectors of each co-occurrence matrix. These  $100 \times 100$  similarity matrices are visualized as square heatmaps where perfect diagonal alignment would indicate identical principal concept axes, while off-diagonal patterns reveal would revolve rotations and mismatches in compositional geometry.

## B Custom RA-SAE Training Details

This section summarizes the training details of our custom archetypal SAE (RA-SAE) [56]. The model is open-sourced at <https://conceptual-blindspots.github.io/>.

**Dataset.** The auto-encoder is trained on the complete ImageNet-1k training split, ( $\approx 1.28$  M RGB images. Each image is converted to 261 visual tokens using DINOv2 [135]; tokens are fed to the SAE without class or position embeddings. The total number of training tokens is therefore  $50 \times 1.28M \times 261 \approx 1.67 \times 10^{10}$ .

**Dictionary.** The dictionary has 32,000 concept dimensions. For the sparse activation rule, top- $k$  masking with  $k = 5$ , is used; activations outside the largest five per input are set to 0. The weights are initialized using Xavier/Glorot. The training is conducted at mixed precision (fp16), with the last ten epochs performed at full precision.

**Optimizer and Schedule.** The model is trained for 50 epochs using base AdamW ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ) optimizer is employed with weight decay set to  $10^{-5}$ . Linear warm-up is applied on the first 5% of steps, followed by cosine decay from  $\eta_{\text{max}} = 5 \times 10^{-4}$  to  $\eta_{\text{final}} = 10^{-6}$ . MSE loss is used alongside an auxiliary term penalizing activations that never enter the top- $k$  set, where  $\lambda = 10^{-5}$ .

## C Computational Resources

This section summarizes the GPU resources used for training and experiments in support of this paper. In total, we used approximately 202 GPU-hours on NVIDIA H100s and H200s.

**RA-SAE.** Trained for approximately 24 GPU-hours on three NVIDIA H100s.

**Synthesized Images.** Generating the full  $D_{g\theta}$  (see Appendix A) took roughly 5 hours per generator when distributed across four NVIDIA H200 GPUs. With five generators, this totaled approximately 100 GPU-hours on a NVIDIA H200.

● **Distribution-Level Analysis.** Extracting energy differences at the distribution level took about 3 hours per generator on a single NVIDIA H200 GPU (total  $\sim 15$  GPU-hours).

● **Datapoint-Level Analysis.** Computing datapoint-level energy differences, ranking concepts per datapoint, and ranking datapoints per concept also took approximately 3 hours per generator on one H200 GPU (total  $\sim 15$  GPU-hours).

## D Autointerpretability

In an effort to automate the interpretation of concepts identified by the RA-SAE, we qualitatively evaluated the ability of Vision-Language Models (VLMs) to describe them. In particular, we queried ChatGPT 4o via the OpenAI API using various prompts and formats of an image  $x''$  in which the given concept was salient, recording the description generated by the VLM. In this section, we highlight the findings from our exploration. For evaluation purposes, we allowed the VLM to provide its full reasoning; to make this a scalable solution, a simple modification of the prompt—asking the VLM to provide only the description—would filter out the reasoning.

First, we provided the VLM with two images: a raw version of  $x''$  and  $x''$  with a heatmap superimposed to localize the intensity of the concept:

Input	Output
 You are given two images: (1) raw image and (2) the same image with a heatmap of areas that are salient for a concept present in the image, as identified by a DINOv2-based SAE. Describe what is this concept.	The concept identified by the DINOv2-based SAE in these images appears to correspond to "ring" or "metal ring" as the salient object. <b>Reasoning:</b> <ul style="list-style-type: none"> <li>- In the original image, a sleek, brushed-metal ring is centered on a wooden surface.</li> <li>- In the heatmap overlay, the regions that are most activated (least transparent with bright color blocks) are tightly focused on the ring itself, particularly the outer surface and inner reflections.</li> <li>- The wooden background is mostly ignored except for minor scattered activations, confirming that the concept is not about the background texture.</li> <li>- The consistent and strong focus on the metallic, circular object suggests that the model recognizes this as the defining element—likely trained to detect the shape, material, and reflectivity typical of rings.</li> </ul>

As seen in the example, the VLM did not interpret the heatmap visualization correctly. Instead of describing the wood texture in the background—the actual salient region—it described the ring in the foreground, which received no attention from the RA-SAE.

To determine whether this was a prompting issue, we tested the same scenario again—this time using an expanded prompt that explained the heatmap in greater detail:

Input	Output
 You are given two images: (1) raw image and (2) the same image with a superimposed heatmap of areas that contain a concept identified by a DINOv2-based SAE. Red areas are the most salient (contain the concept); blue areas are the least salient (do not contain the concept). Describe what is this concept.	The concept identified by the DINOv2-based SAE in these images appears to be "ring" or more specifically "metal band/jewelry ring." <b>Here's the reasoning:</b> <ul style="list-style-type: none"> <li>- In the heatmap image, the red and yellow (most salient) regions align very closely with the circular metal ring at the center of the image.</li> <li>- The ring is the only distinct object in the photo and is clearly segmented from the background (wooden surface), which is mostly marked in blue and green (least salient).</li> <li>- The attention does not focus on the texture of the background or lighting but is highly localized to the curved, metallic object.</li> </ul> Thus, the concept is almost certainly the metallic ring—its shape, reflectivity, and material seem to trigger the sparse autoencoder's activation, making it a learned concept in the SAE dictionary.

This prompt modification did not resolve the issue, as the VLM still identified the ring as the prominent concept in the image. Several other prompting techniques, including those describing the heatmap in even greater depth (e.g., by including the code used for its generation), were explored with similar results.

We next tested providing the VLM only with the heatmap-superimposed  $x''$  and a similar prompt explaining how to interpret the heatmap:

Input	Output
 <div style="background-color: #f0e6d2; padding: 5px; margin-top: 10px;"> <p>You are given an image with a superimposed heatmap of areas that contain a concept identified by a DINOv2-based SAE. Red areas are the most salient (contain the concept); blue areas are the least salient (do not contain the concept). Describe what is this concept.</p> </div>	<div style="background-color: #e0e0e0; padding: 10px;"> <p>The concept identified by the DINOv2-based SAE in this image appears to be "<b>ring</b>" or more broadly, "<b>circular metallic object</b>". The most salient (red) areas are focused tightly around the ring itself—especially along its edges and reflective interior—indicating that the model has strongly associated those regions with the concept.</p> <p>The surrounding wooden surface is mostly in blue and green, which shows it is not relevant to the concept, supporting the idea that the model is isolating the ring as the key object of interest.</p> </div>

Once again, the result failed to identify the wooden background as the prominent concept in the image and instead focused on the ring. As in the previous case, we experimented with various prompting strategies using this image format, but the results remained largely unchanged.

From these observations, we concluded that the heatmap alone is insufficient to guide the VLM’s attention to the salient areas. The presence of less relevant regions appears to distract the model. To address this, we replaced the heatmap-superimposed  $x''$  with an alpha-masked version of  $x''$ :

Input	Output
 <div style="background-color: #f0e6d2; padding: 5px; margin-top: 10px;"> <p>You are given an image where only areas that contain a concept identified by a DINOv2-based SAE are preserved; areas without the concept were removed. Describe what is this concept.</p> </div>	<div style="background-color: #e0e0e0; padding: 10px;"> <p>The concept identified and preserved in this image appears to be "<b>wood texture</b>" or "<b>wood grain</b>." The retained areas show consistent patterns typical of wooden surfaces—linear, fibrous streaks in earthy tones like beige, gray, and brown. The selective masking highlights parts of the image where this natural texture is most prominent, suggesting the SAE learned to isolate the visual features of wood (likely flooring or a table).</p> </div>

As shown in the example, the VLM was now able to correctly identify `wood texture` as the dominant concept. We observed consistent success across a range of concepts and images.

The formatting of  $x''$  and the prompt shown above yielded the most reliable results in our qualitative evaluation. However, we note that this evaluation is limited by its qualitative nature (due to the absence of ground truth annotations) and its focus on a single VLM. We hope future work on the autointerpretability of SAE concepts can build on and expand this analysis.

## E Additional Results: Model-Specific Blindspots

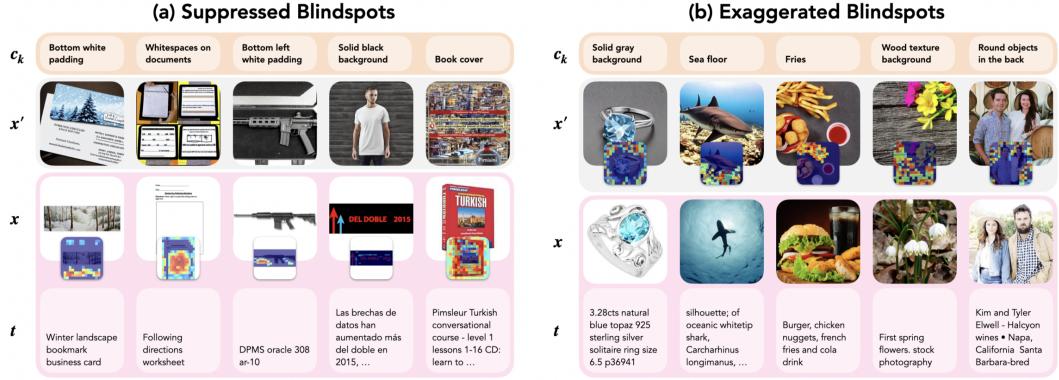


Figure 14: Examples of conceptual blindspots in **Stable Diffusion 1.5**. For each concept, the prototypical natural (for suppressed blindspots) or synthesized (for exaggerated blindspots), based on the highest absolute activation, is shown. The spatial heatmap for the concept is superimposed atop the image.

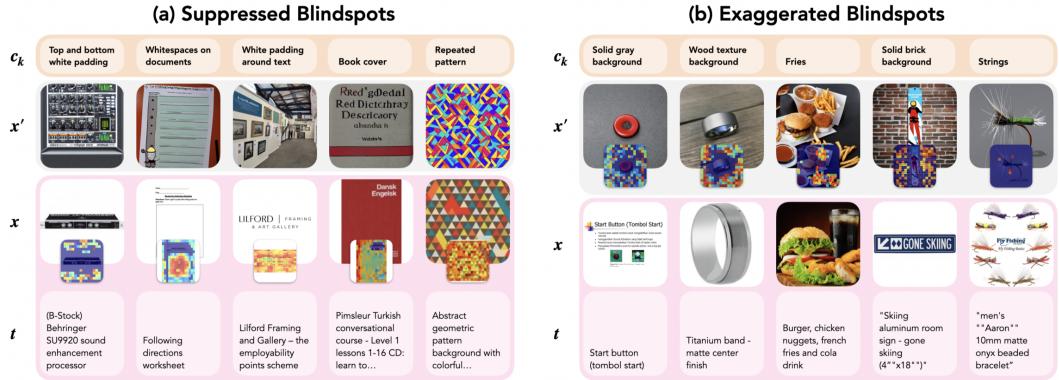


Figure 15: Examples of conceptual blindspots in **Stable Diffusion 2.1**. For each concept, the prototypical natural (for suppressed blindspots) or synthesized (for exaggerated blindspots), based on the highest absolute activation, is shown. The spatial heatmap for the concept is superimposed atop the image.

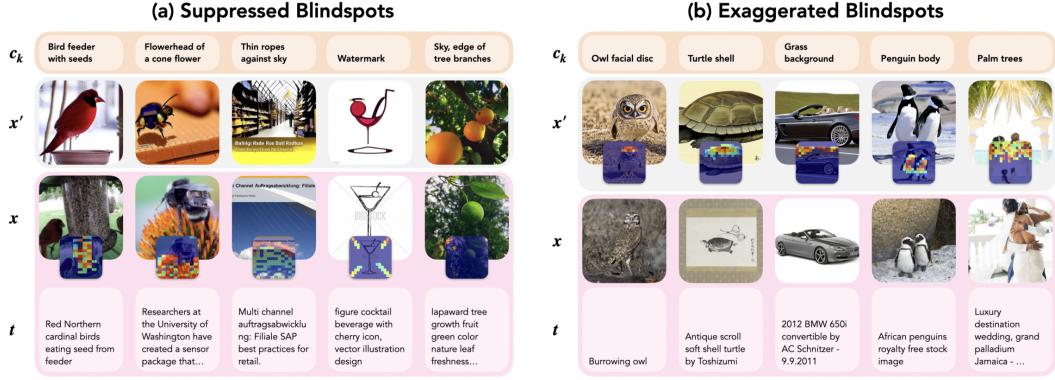


Figure 16: Examples of conceptual blindspots in **Kandinsky**. For each concept, the prototypical natural (for suppressed blindspots) or synthesized (for exaggerated blindspots), based on the highest absolute activation, is shown. The spatial heatmap for the concept is superimposed atop the image.

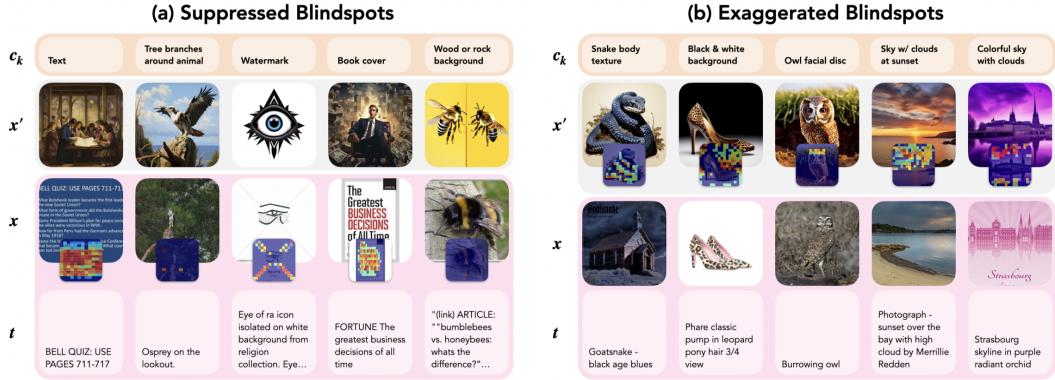


Figure 17: Examples of conceptual blindspots in **PixArt**. For each concept, the prototypical natural (for suppressed blindspots) or synthesized (for exaggerated blindspots), based on the highest absolute activation, is shown. The spatial heatmap for the concept is superimposed atop the image.

## F Additional Results: Higher-order Blindspots with Compositional Discrepancy

Thus far, our analysis has centered on individual concept activations. Yet visual scenes are rarely composed of isolated concepts; instead, they are structured through rich and structured co-occurrence patterns that encode compositional semantics. We now examine whether generative models capture this higher-order structure by analyzing the co-activation matrix  $Z^T Z$ , which reflects pairwise correlations in concept usage.

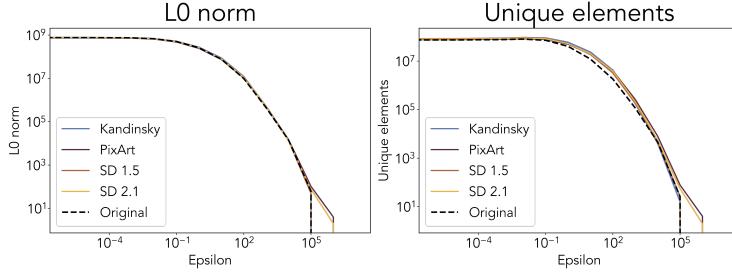


Figure 18: ● **Sparsity and Structural Divergence.** On the left:  $\ell_0$  norm of the co-occurrence matrix  $ZZ^T$  as a function of  $\epsilon$  (threshold), indicating how many entries remain active in each model. On the right: Number of unique entries in the synthesized distribution relative to the natural distribution. All evaluated models preserve global sparsity structure, but diverge in activation content.

Surprisingly, when assessed at the level of binary structure, diffusion models approximate the global sparsity of the natural co-occurrence matrix with high fidelity. As shown in Fig. 18 (left), the  $\ell_0$  norm of  $Z^T Z$  – thresholded at varying  $\epsilon$  values – tracks closely between the natural and synthesized distributions across all models. This indicates that the gross connectivity of the conceptual graph, i.e., which concepts tend to co-activate at all, is well preserved. Formally, one can deem  $Z^T Z$  as the adjacency matrix of a weighted, undirected graph over concepts, where edge weights reflect co-activation strength across the dataset.

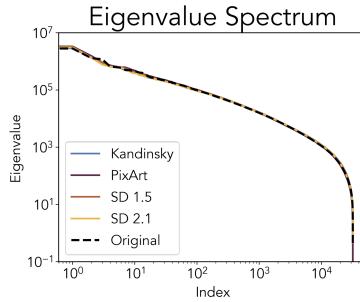


Figure 19: ● **Spectral Structure of Co-occurrence.** Log-log plot of the eigenvalue spectra from co-occurrence matrices  $ZZ^T$  across models and the natural distribution. All evaluated models match the heavy-tailed decay of the natural distribution.

However, as illustrated in Fig. 18 (right), the specific content of these co-activations diverges: a substantial portion of entries in the model-generated  $Z^T Z$  are not shared with the natural baseline. This suggests that while the capacity for compositionality is retained, the identity of active pairings may shift, potentially reflecting model specific inductive biases or training artifacts. To probe the internal structure of these co-occurrence patterns, we turn to spectral analysis. Fig. 19 shows the eigenvalue spectra of the co-occurrence matrices for each model and the natural distribution. All spectra exhibit a heavy-tailed decay, consistent with power-law behavior, indicating that generative models preserve the overall rank structure and variance allocation across conceptual dimensions.

Further, we examine the alignment of dominant conceptual directions via cosine similarity heatmaps between the top 100 eigenvectors of the synthesized and natural co-occurrence matrices (Fig. 20). While all models exhibit partial diagonal alignment—implying overlap in principal concept axes—the

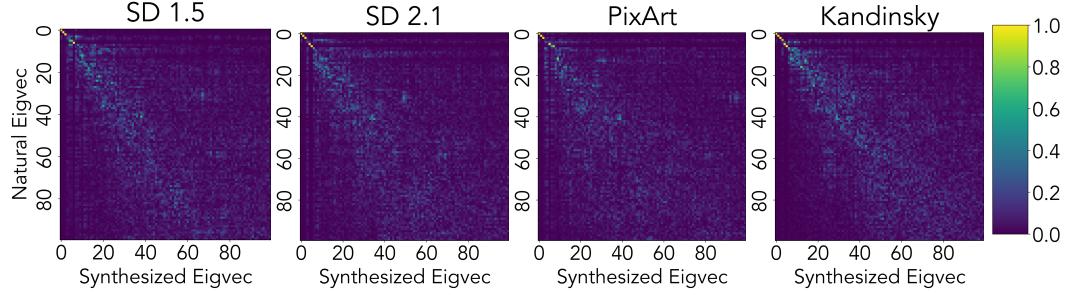


Figure 20: ● **Concept Basis Similarity.** Cosine similarity heatmaps between the top 100 eigenvectors of the natural and synthesized co-occurrence matrices  $\bar{Z}\bar{Z}^T$ . Diagonal structure shows alignment of dominant conceptual directors, with varying degrees of alignment across the four models.

off-diagonal entries reveal rotations and mismatches in higher modes, reflecting evident deviations in compositional geometry.

Together, these findings reveal that diffusion models approximate the global shape of concept co-activation surprisingly well, yet deviate in subtle and structured ways when examined through the spectral lens. Such higher-order discrepancies may underpin failures in generating coherent, multi-object scenes or relational concepts.

## G Additional Results: Blindsight Stress Testing

To stress-test the blindspots identified by our method, we gathered a range of prompts describing these blindspots and used them to generate many images. We then contrasted the outputs from models in which the concept was identified as a blindspot with those in which it was not.

Specifically, ChatGPT-4o was prompted as follows: *I want to generate an image of the following concept: "<blindspot>". Suggest 50 prompts highlighting this concept to be used as input for a text-to-image model. Return these as a list of strings in Python.* Five images were generated per prompt and analyzed using our custom RA-SAE model (see Appendix B), which ranked them by the intensity with which the desired concept appeared. All images were then manually reviewed to determine whether the blindspot was successfully depicted.

As seen in the following examples of suppressed concepts, while some aspects of the target concept occasionally appeared (e.g., a holder or string for the `bird feeder` blindspot and a round hole for the `glossy DVD disc` blindspot), the models generally failed to generate the full concept. This aligns with our method’s assessment and supports the validity of the stress test.

### G.1 Bird Feeder Blindsight in Kandinsky



Figure 21: Examples of images generated with various prompts involving the `bird feeder` concept as a part of the stress testing. In Kandinsky, our method identified this concept as a suppressed conceptual blindspot, which matches the observed behavior: the model is unable to generate a corresponding image. By contrast, SD 1.5, in which this concept was not identified as a blindspot, is able to generate this concept.

## G.2 Glossy DVD Disc Blindspot in SD 1.5

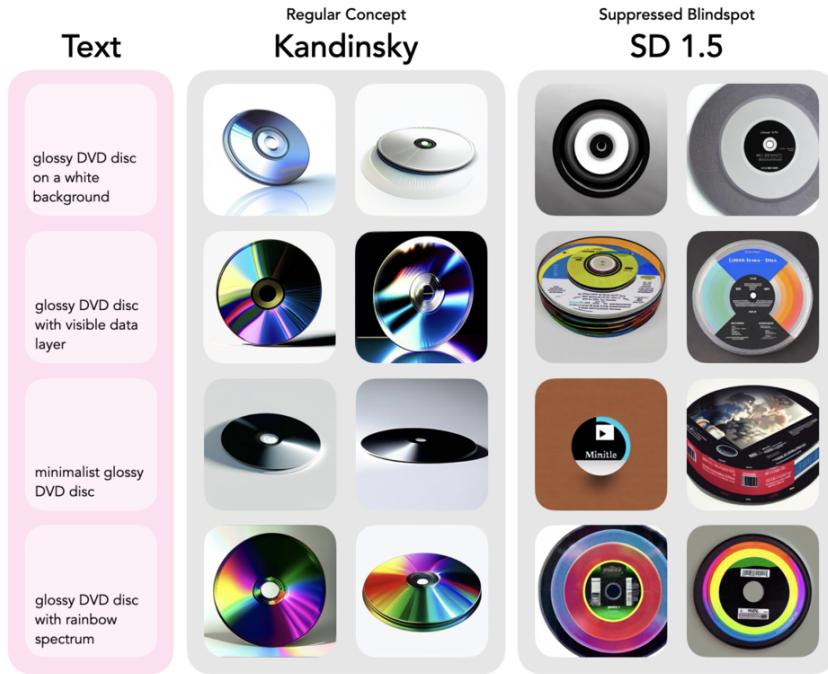


Figure 22: Examples of images generated with various prompts involving the glossy DVD disc concept as a part of the stress testing. In SD 1.5, our method identified this concept as a suppressed conceptual blindspot, which matches the observed behavior: the model is unable to generate a corresponding image. By contrast, Kandinsky, in which this concept was not identified as a blindspot, is able to generate this concept.

## H Additional Examples: Synthesized Images

Text	Natural	SD 1.5	SD 2.1	PixArt	Kandinsky
New Modern Fiberglass Plastic Swimming Pool Entrance Step With Clean Fresh Refreshing Blue Water On					
Kylie - Yorkie - Rolly Teacup Puppies					
gwen stefani wind it up mp3:					
office bundle package					
Stagecoach Springs Road, Custer image					
Christ On The Cross Acrylic Print by Delacroix					
Picture of Outdoor Swing Chair / Bed Canopy Patterned Arch - Coffee					
Dogs and Their People					
carbonara spaghetti on wood background					

Figure 23: Additional image-caption pair examples from LAION-5B with matching images generated with the same prompt by SD 1.5, SD 2.1, PixArt, and Kandinsky.

## I Concentration bounds for $\delta$

In our experiments, we estimate  $\delta(k)$  using  $n = 10,000$  paired samples for each concept. While this budget is modest, it raises the natural question of whether it suffices to obtain reliable estimates. To address this, we derive a concentration bound on the empirical estimator  $\widehat{\delta}_n(k)$  using McDiarmid's inequality [142]. The resulting bound is tight and demonstrates that even with relatively few samples, we can obtain fast and accurate estimates of concept bias.

**Theorem 4** (Concentration of  $\widehat{\delta}_n(k)$ ). *We assume that the concept score  $\xi_k(\mathbf{x})$  takes values in  $[a, b]$  almost surely for all images  $\mathbf{x}$  drawn from either  $\mathcal{D}_{\mathcal{X}}$  or  $\mathcal{D}'_{\mathcal{X}}$ . Let  $n$  paired samples  $(\mathbf{x}_i, \mathbf{x}'_i)_{i=1}^n$  be drawn independently with  $\mathbf{x}_i \sim \mathcal{D}_{\mathcal{X}}$  and  $\mathbf{x}'_i \sim \mathcal{D}'_{\mathcal{X}}$ , and define the empirical estimator*

$$\widehat{\delta}_n(k) := \sigma \left( \frac{1}{n} \sum_{i=1}^n \xi_k(\mathbf{x}_i) - \frac{1}{n} \sum_{i=1}^n \xi_k(\mathbf{x}'_i) \right).$$

Let  $M := b - a$  and  $L := M/4$ . Then for every  $\varepsilon > 0$ , the deviation satisfies

$$\mathbb{P} \left( \left| \widehat{\delta}_n(k) - \delta(k) \right| > \varepsilon \right) \leq 2 \exp \left( -\frac{2n\varepsilon^2}{L^2} \right).$$

*Proof.* The function  $x \mapsto \sigma(x)$  is 1/4-Lipschitz, since  $|\sigma'(x)| \leq 1/4$  for all  $x$ .

Viewing  $\widehat{\delta}_n(k)$  as a function of the  $2n$  independent variables  $(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{x}'_1, \dots, \mathbf{x}'_n)$ , changing a single argument alters the inner difference of means by at most  $M/n$ , and the outer sigmoid scales this by at most 1/4. Hence, the bounded difference constant for each coordinate is  $(M/n)(1/4) = L/n$ .

By McDiarmid's inequality [142],

$$\mathbb{P} \left( \left| \widehat{\delta}_n(k) - \delta(k) \right| > \varepsilon \right) \leq 2 \exp \left( -\frac{2\varepsilon^2}{\sum_{j=1}^{2n} (L/n)^2} \right) = 2 \exp \left( -\frac{2n\varepsilon^2}{L^2} \right),$$

which proves the claim.  $\square$

Practically, most concept scores  $\delta(k)$  are sparse, with the majority concentrated near zero and only a few reaching values up to 10. The concentration bound shows that even for the largest observed biases, a sample size of  $n = 10,000$  yields estimates of  $\widehat{\delta}_n(k)$  that deviate from the true value by no more than a small  $\varepsilon$  with high probability. This justifies our sampling strategy and confirms that accurate bias measurements are attainable with limited data.

## J Monotonicity and Calibration-Free Interpretation of $\delta$

Our goal when analyzing blind spots is to rank concepts by the severity of their generative bias. In practice, we use the score  $\delta(k)$  for this purpose. However, one may wonder whether such a score introduces distortions relative to more direct quantities such as the energy gap or the odds ratio. The following result establishes that  $\delta(k)$  is a strictly increasing reparameterization of both, and therefore inherits their ordering. This guarantees that no calibration is needed when using  $\delta(k)$  to rank concepts.

**Theorem 5** (Monotonicity and Calibration of  $\delta_{g_\theta \leftrightarrow G}$ ). *For every concept index  $k$  define the energy gap*

$$\Delta_k = \mathbb{E}_{\mathbf{x}' \sim \mathcal{D}'_X} [\xi_k(\mathbf{x}')] - \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_X} [\xi_k(\mathbf{x})],$$

*the associated odds ratio  $\rho_k = \exp(\Delta_k)$ , and the energy-difference score*

$$\delta(k) = \frac{1}{1 + \exp(-\Delta_k)} = \frac{\rho_k}{1 + \rho_k}.$$

*Then  $\delta(k)$  is a strictly increasing bijection of both  $\Delta_k$  and  $\rho_k$ , so ranking concepts by any one of  $\delta(k)$ ,  $\Delta_k$ , or  $\rho_k$  produces exactly the same ordering.*

*Proof.* The logistic sigmoid satisfies  $\sigma'(x) = \sigma(x)(1 - \sigma(x)) > 0, \forall x \in \mathbb{R}$ ; hence  $\sigma$  and therefore  $\delta(k) = \sigma(\Delta_k)$  grow strictly with  $\Delta_k$ . Because the exponential map is also strictly increasing and bijective  $\mathbb{R} \rightarrow (0, \infty)$ , setting  $\rho_k = \exp(\Delta_k)$  preserves order and gives  $\Delta_k = \log \rho_k$ . Substituting this identity into  $\sigma$  yields  $\delta(k) = \sigma(\log \rho_k) = \rho_k / (1 + \rho_k)$ , which is the composition of two strictly increasing bijections and is therefore itself strictly increasing and bijective in  $\rho_k$ . Since strict monotonic functions never reverse inequalities, the three quantities share the same total order over concepts.  $\square$

Thus, ranking concepts by  $\delta(\cdot)$  is strictly equivalent to ranking them by energy gap or by conceptual generation odds  $\rho_k$ . No calibration is necessary, and all three quantities preserve the same total ordering over concepts.

## K Stability of FID Under SAE Embeddings

In this section we establish a quantitative relationship between the Fréchet Inception Distance (FID) computed in the original activation space (of dimension  $d$ ) and the FID after applying a (potentially overcomplete) SAE dictionary  $\mathbf{D} \in \mathbb{R}^{k \times d}$  with  $k \gg d$ . Throughout we assume that  $\mathbf{D}$  has orthonormal columns but is not necessarily square, i.e.

$$\mathbf{D}^\top \mathbf{D} = \mathbf{I}_d, \quad \text{while} \quad \mathbf{D}\mathbf{D}^\top \neq \mathbf{I}_k.$$

We start by recalling a simple fact: if  $\mathbf{D}$  is not overcomplete, orthogonal and  $k = d$ , then we have an isometry between  $\mathbf{A}$  and  $\mathbf{Z}$ , implying that the FID is perfectly preserved. However, this case is not realistic, we will then turn the overcomplete case, and show we can bound FID by the extreme singular value of  $\mathbf{D}$ . We will work with the Wasserstein-2 metric  $\mathcal{W}_2$ , noting that FID is just  $\mathcal{W}_2^2$  specialised to Gaussians.

For a probability measure  $\mu$  on  $\mathbb{R}^d$  we write  $\mathbf{D}_\#\mu$  for its push-forward under  $\mathbf{D}$ , i.e.  $\mathbf{D}_\#\mu(\mathbf{z}) = \mu(\mathbf{D}^{-1}\mathbf{z})$ . Denote by  $\sigma_{\min}$  and  $\sigma_{\max}$  the minimal and maximal singular values of  $\mathbf{D}$ , equivalently the square-roots of the extremal eigenvalues of  $\mathbf{D}\mathbf{D}^\top$ :

$$\sigma_{\min}^2 \mathbf{I}_k \preceq \mathbf{D}\mathbf{D}^\top \preceq \sigma_{\max}^2 \mathbf{I}_k.$$

Empirically one usually finds  $\sigma_{\min}, \sigma_{\max} \approx 1$ , but the proof does not rely on that. We will start by a simple lemma in the case where  $\mathbf{D}$  is not overcomplete.

**Lemma 1** (Isometry under exact orthogonality). *Suppose  $k = d$  and  $\mathbf{D}^\top \mathbf{D} = \mathbf{D}\mathbf{D}^\top = \mathbf{I}_d$ . Then  $\mathbf{D}$  is an isometry:  $\|\mathbf{D}\mathbf{v}\|_2 = \|\mathbf{v}\|_2$  for all  $\mathbf{v} \in \mathbb{R}^d$ . Consequently, for any probability measures  $\mu, \nu$  on  $\mathbb{R}^d$  with finite second moment,*

$$\mathcal{W}_2(\mathbf{D}_\#\mu, \mathbf{D}_\#\nu) = \mathcal{W}_2(\mu, \nu).$$

*Proof.* Orthogonality of  $\mathbf{D}$  implies preservation of the Euclidean norm, and push-forward commutes with the map inside the  $\mathcal{W}_2$  infimum; the integrand is unchanged, so the infimum value is identical.  $\square$

This case, however, is quite unrealistic as SAE usually rely on the overcompleteness to extract meaningful and interpretable concepts. In the overcomplete case,  $\mathbf{D}$  is no longer orthonormal, but we can still have column-orthonormal dictionary. We will use that to show that we can bound using the extremal singular value of  $\mathbf{D}^\top \mathbf{D}$ .

**Theorem 6** (FID under column orthogonal embeddings). *Let  $\mathbf{D} \in \mathbb{R}^{k \times d}$  satisfy  $\mathbf{D}^\top \mathbf{D} = \mathbf{I}_d$  and denote by  $0 < \sigma_{\min} \leq \sigma_{\max}$  the extreme singular values of  $\mathbf{D}\mathbf{D}^\top$ . Given two data matrices  $\mathbf{A}, \mathbf{A}' \in \mathbb{R}^{n \times d}$  (rows are sample vectors), set  $\mathbf{Z} = \mathbf{A}\mathbf{D}^\top \in \mathbb{R}^{n \times k}$  and  $\mathbf{Z}' = \mathbf{A}'\mathbf{D}^\top \in \mathbb{R}^{n \times k}$ . Then*

$$\sigma_{\min}^2 \text{FID}(\mathbf{A}, \mathbf{A}') \leq \text{FID}(\mathbf{Z}, \mathbf{Z}') \leq \sigma_{\max}^2 \text{FID}(\mathbf{A}, \mathbf{A}').$$

*Proof.* Write  $\mu$  for the empirical measure of  $\mathbf{A}$  and  $\nu$  for that of  $\mathbf{A}'$ , i.e.

$$\mu = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{A}_{i,:}}, \quad \nu = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{A}'_{i,:}}.$$

For any coupling  $\pi \in \Pi(\mu, \nu)$  (i.e. a probability measure on  $\mathbb{R}^d \times \mathbb{R}^d$  with marginals  $\mu, \nu$ ) we have, by the extremal singular value bound,

$$\sigma_{\min}^2 \|\mathbf{x} - \mathbf{y}\|_2^2 \leq \|\mathbf{D}(\mathbf{x} - \mathbf{y})\|_2^2 \leq \sigma_{\max}^2 \|\mathbf{x} - \mathbf{y}\|_2^2, \quad \forall (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^d \times \mathbb{R}^d.$$

Integrating with respect to an arbitrary coupling  $\pi \in \Pi(\mu, \nu)$  yields

$$\sigma_{\min}^2 \int \|\mathbf{x} - \mathbf{y}\|_2^2 d\pi \leq \int \|\mathbf{D}(\mathbf{x} - \mathbf{y})\|_2^2 d\pi \leq \sigma_{\max}^2 \int \|\mathbf{x} - \mathbf{y}\|_2^2 d\pi.$$

The middle integral is exactly the transport cost of the pushed-forward coupling  $(\mathbf{D} \times \mathbf{D})_\#\pi$  between  $\mu_D := \mathbf{D}_\#\mu$  and  $\nu_D := \mathbf{D}_\#\nu$ . Because the inequalities hold for every  $\pi$ , they hold in particular

for the optimal couplings attaining  $\mathcal{W}_2(\boldsymbol{\mu}, \boldsymbol{\nu})$  and  $\mathcal{W}_2(\boldsymbol{\mu}_D, \boldsymbol{\nu}_D)$ , though these two optima need not coincide. Taking the infimum over  $\boldsymbol{\pi}$  term-wise makes this explicit:

$$\sigma_{\min}^2 \inf_{\boldsymbol{\pi} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu})} \int \|x - y\|_2^2 d\boldsymbol{\pi} \leq \inf_{\boldsymbol{\pi} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu})} \int \|D(x - y)\|_2^2 d\boldsymbol{\pi} \leq \sigma_{\max}^2 \inf_{\boldsymbol{\pi} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu})} \int \|x - y\|_2^2 d\boldsymbol{\pi}.$$

Hence

$$\sigma_{\min}^2 \mathcal{W}_2^2(\boldsymbol{\mu}, \boldsymbol{\nu}) \leq \mathcal{W}_2^2(\boldsymbol{\mu}_D, \boldsymbol{\nu}_D) \leq \sigma_{\max}^2 \mathcal{W}_2^2(\boldsymbol{\mu}, \boldsymbol{\nu}).$$

Recognising  $\text{FID}(\cdot, \cdot) = \mathcal{W}_2^2(\cdot, \cdot)$  for the Gaussian surrogate and plugging in  $(\mathbf{A}, \mathbf{A}')$  (resp.  $(\mathbf{Z}, \mathbf{Z}')$ ) finishes the proof.  $\square$

Essentially, theorem 6 tells us that applying a column-orthogonal overcomplete SAE dictionary cannot distort Fréchet Inception Distance by more than the square of its extremal singular values. When  $D$  is nearly orthogonal – empirically we usually found that  $\sigma_{\min}, \sigma_{\max} \approx 1$  – the result implies that FID measured in the SAE feature space is essentially close to the canonical FID.  $\square$