

CANONICAL LATENT REPRESENTATIONS IN CONDITIONAL DIFFUSION MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Conditional diffusion models (CDMs) have shown impressive performance across a range of generative tasks. Their ability to model the full data distribution has opened new avenues for analysis-by-synthesis in downstream discriminative learning. However, this same modeling capacity causes CDMs to entangle the class-defining features with irrelevant context, posing challenges to extracting robust and interpretable representations. To this end, we introduce *Canonical Latent Representation Identifier* (CLARID), a training-free procedure to identify *Canonical Latent Representations* (CanoReps), latent codes whose internal CDM features preserve essential categorical information while discarding non-discriminative signals. When decoded, CanoReps produce representative samples for each class, offering an interpretable and compact summary of the core class semantics with minimal irrelevant details. Exploiting CanoReps, we develop a novel diffusion-based feature distillation paradigm, *CaDistill*. While the student has full access to the training set, the CDM as teacher transfers core class knowledge only via CanoReps, which amounts to merely 10% of the training data in size. After training, the student achieves strong adversarial robustness and generalization ability, focusing more on the class signals instead of spurious background cues. Our findings suggest that CDMs can serve not just as image generators but also as compact, interpretable teachers that can drive robust representation learning.

1 INTRODUCTION

Diffusion models (DMs) excel at generative modeling of images (Dhariwal & Nichol, 2021; Ho et al., 2020; Rombach et al., 2022; Song et al., 2021a;b). When conditioned on class labels (Bao et al., 2023; Ho & Salimans; Peebles & Xie, 2023) or text prompts (Podell et al., 2024; Rombach et al., 2022), conditional diffusion models (CDMs) faithfully generate samples with desired characteristics of the condition. This generative capability has sparked a wave of analysis-by-synthesis approaches (Baranchuk et al., 2022; Chen et al., 2024b; Li et al., 2023b; Meng et al., 2024; Mukhopadhyay et al., 2024; Xiang et al., 2023; Yang & Wang, 2023; Zhang et al., 2024; Zhao et al., 2023b), where DMs are used to probe or improve downstream discriminative tasks. However, a key challenge remains: since DMs model the full data distribution, they often encode redundant or irrelevant information, which can obscure the discriminative signal. For example, in the *Tench* row of Figure 1, modifying the CDM latent code of the training sample changes the angler in the background while the fish remains almost unchanged, showing that the model encodes background cues that correlate with the class but not semantically essential to it. This entanglement between class semantics and extraneous factors limits interpretability and hinders the effective use of CDMs in representation learning. This motivates our central question:

How can we identify the underlying core categorical semantics in a conditional diffusion model?

We answer this by introducing **Canonical LATent Representation IDentifier** (CLARID), a method for identifying the latent codes in CDMs that capture essential categorical information and filter out irrelevant details. We call the resulting latent codes **Canonical Latent Representations** (CanoReps). We begin with the assumption that the essential semantics of a class typically lie on a low-dimensional manifold embedded in the high-dimensional data space (Collins et al., 2018; Wang et al.; Zhang et al., 2019). We postulate that such semantic manifolds also exist within the latent space of diffusion models. Our key insight is that altering the latent code along the tangent directions of

the manifold, referred to as extraneous directions, modifies visual appearance without affecting class identity. We find that projecting out the extraneous directions in the latent space of CDMs effectively eliminates class-irrelevant factors such as background clutter or co-occurring objects from other categories. When projected back to the data space, CanoReps produce representative samples of each category, namely *Canonical Samples*, providing an intuitive and interpretable summary of the essential categorical semantics. Additionally, the internal CDM features of CanoReps, *i.e.* *Canonical Features*, contain mostly the core class information. We first validate our method in a toy hierarchical generative model, where CLARID recovers a low-dimensional class manifold while standard classifier-free guidance (CFG) produces dispersed, high-likelihood samples. Scaling up, we apply CLARID to ImageNet-pretrained CDMs and develop strategies for finding the projection time step and the number of extraneous directions. Our method also generalizes to text-conditioned DMs and is compatible with different diffusion samplers.

Building on the discovery of CanoRep, we propose a novel feature distillation paradigm, *CaDistill*. This method leverages the interpretable nature of CanoReps, which encapsulates the core semantics of each class, to supervise a student network. *CaDistill* aligns the student network’s representations on both Canonical Samples and original training samples with Canonical Features using a novel feature distillation loss, which helps the student network encode the core class information. The student network’s representations of the original training samples are also forced to be close to those of the Canonical Samples, treating them as anchors in the student’s representation space. The student learns on the full training set, while the teacher CDM transfers the essential class knowledge by only exploiting CanoReps, which amounts to merely 10% of the training data in size. In contrast, existing state-of-the-art methods require transferring the teacher CDM’s knowledge using the entire training set. *CaDistill* improves the student’s adversarial robustness as well as the generalization capability on CIFAR10 (Krizhevsky et al., 2009) and ImageNet (Deng et al., 2009). Our contributions are as follows:

- We introduce **Canonical Latent Representation** (CanoRep) in CDMs—latent codes whose internal CDM features, *i.e.* *Canonical Features*, encapsulate core categorical semantics with minimal irrelevant signals. When decoded to the data space, these latent codes produce *Canonical Samples*, which serve as compact and interpretable prototypes for each class.
- To extract *CanoReps*, we propose **Canonical Latent Representation Identifier** (CLARID), a method that identifies these latent codes by projecting out non-discriminative directions in the CDM’s latent space. The optimal configurations of CLARID are selected through a systematic analysis of the CDM’s features.
- Leveraging the CanoReps, we develop a novel diffusion-based feature distillation paradigm, *CaDistill*. While the student is being trained on the full training set, the CDM as the teacher transfers essential class knowledge only via exploiting CanoReps, which amounts to merely 10 % of the data. The resulting student network achieves strong adversarial robustness and generalization performance, focusing more on the core class information.

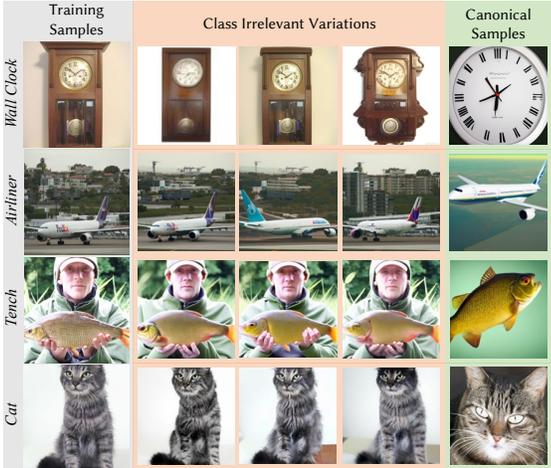


Figure 1: Conditional diffusion models (CDMs) encode both **core class features** and **irrelevant context** in a CDM. **Left**: training samples. **Middle**: samples obtained by modifying CDM latent codes that preserve the class label but alter class-irrelevant parts. **Right**: Canonical Samples produced by CLARID, which retain the **class-defining content** while removing **extraneous context**. CLARID benefits the extraction of robust and interpretable representations from the CDMs.

2 RELATED WORKS

2.1 INTERPRETABILITY IN DIFFUSION MODELS

Recent research in diffusion models (DMs) reveals the interpretable semantic information in them. We categorize these efforts into two main groups. The first line of this work focuses on semantic editing by manipulating the reverse diffusion trajectory to produce semantically meaningful changes in generated images (Chen et al., 2024a; Haas et al., 2024; Kwon et al., 2023; Park et al., 2023a). Kwon et al. (2023) uncover a semantic space inside a pre-trained DM, termed h -space, where particular vector directions yield high-quality image editing results. Park et al. (2023a) analyze the latent input space, namely x -space, of DMs from a Riemannian geometry perspective. They define the pullback metric on x -space from the h -space Euclidean metric, obtaining certain vector directions in x -space that can yield semantic editing results. Chen et al. (2024a) provides more theoretical insights into this framework and extends it to local editing scenarios.

Another line of work leverages the attention mechanism (Vaswani et al., 2017) in DMs to interpret the conditional information (Chefer et al., 2023; Hertz et al.; Kim et al., 2023; Liu et al., 2024; Xu et al., 2023; Zhao et al., 2023a). Hertz et al. propose that the cross-attention map between the text prompts and image tokens in text-conditioned DMs encodes rich spatial cues. Building on this insight, subsequent studies analyze these attention maps to improve the precision and controllability of semantic image editing (Chefer et al., 2023; Kim et al., 2023; Liu et al., 2024). Other efforts investigate this property in visual recognition tasks such as semantic segmentation (Xu et al., 2023; Zhao et al., 2023a), using the attention maps to interpret the model’s spatial reasoning. All existing methods either focus on image editing or are tailored to a specific DM architecture. We take the first step to uncover the underlying core class semantics in DMs without any supervision. Our method is compatible with different DM architectures and samplers.

2.2 DMs AS TEACHERS IN ANALYSIS-BY-SYNTHESIS

DM-based feature distillation. Recent works show that the intermediate features in DMs contain rich discriminative information (Baranchuk et al., 2022; Chen et al., 2024b; Li et al., 2023b; Meng et al., 2024; Mukhopadhyay et al., 2024; Xiang et al., 2023; Yang & Wang, 2023; Zhang et al., 2024; Zhao et al., 2023b). Here, we focus on the utilization of DMs as teachers in feature distillation frameworks. Li et al. (2023b) proposes a framework in which the intermediate features of the student network are aligned to those of a generative teacher, improving the student’s performance on dense prediction tasks. Yang & Wang (2023) use reinforcement learning to select a proper diffusion timestep for distillation, enhancing the student’s performance in image classification, semantic segmentation, and landmark detection benchmarks.

DMs for data generation and augmentation. DMs faithfully model the full training data distribution, which allows the generation of new training samples or augmentation of the existing ones (Azizi et al.; Bansal & Grover; Fu et al., 2024; Goyal et al., 2021; He et al., 2023; Islam et al., 2024; Sariyildiz et al., 2023; Sehwag et al., 2022; Shama et al., 2024; Trabucco et al., 2024; Wang et al., 2023). Goyal et al. (2021) and Sehwag et al. (2022) improve the robustness of adversarially-trained classifiers by using diffusion-generated data. Bansal & Grover demonstrate that supplementing training data with diffusion-generated images leads to consistent gains on out-of-distribution test sets. Diffusion-generated data is also effective in data-scarce settings, *i.e.*, zero-shot and few-shot learning (Fu et al., 2024; He et al., 2023; Trabucco et al., 2024). Regarding data augmentation, Islam et al. (2024) propose blending images while preserving their labels using pre-trained text-to-image DMs (Rombach et al., 2022). Shama et al. (2024) utilize denoised samples for augmenting the training data, improving the generalization of downstream classifiers.

Despite the progress, current methods use raw diffusion features and outputs, which contain class-irrelevant information. Such irrelevant signals prevent the student from efficiently and accurately learning the class semantics, leading to vulnerable models. On the contrary, our method transfers the core class semantics using CanoReps to the student, enhancing its adversarial robustness and generalization capability. Notably, CanoReps amounts to only 10% of the original data in size.

3 METHOD

An overview of Canonical Latent Representation Identifier (CLARID) is shown in Figure 2. We describe the procedure of finding the CanoRep of a given sample in Section 3.2. We then provide the intuition of the effect of CLARID in a proof-of-concept experiment in Section 3.3.

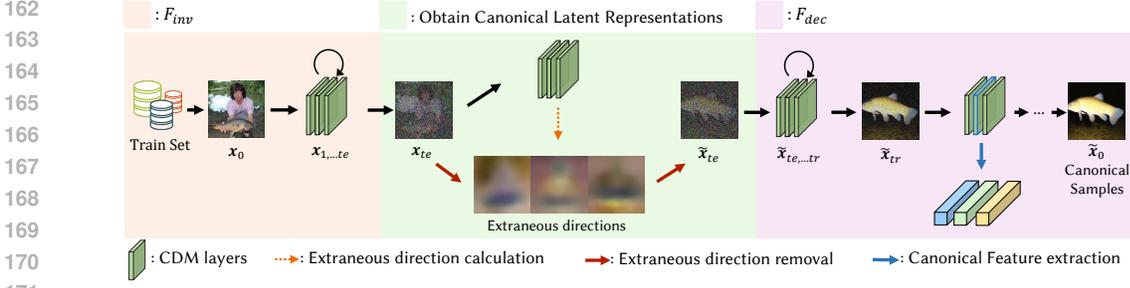


Figure 2: Overview of Canonical Latent Representation Identifier (CLARID). Starting from a training sample x_0 , we invert it using a CDM until $t = t_e$. We compute the extraneous directions using the method described in Section 3.2, and remove them from the inverted sample x_{t_e} to obtain a CanoRep \tilde{x}_{t_e} . We then generate the Canonical Sample \tilde{x}_0 and extract the Canonical Feature at timestep $t = t_r$. The CDM receives the ground truth condition of x_0 in the whole process, here the *Tench* class.

3.1 PRELIMINARIES

Diffusion models (DMs) generate images by learning to reverse a fixed forward diffusion process (Ho et al., 2020). Let x_0 be a training sample and $t \in \{1, \dots, T\}$ denote the diffusion time steps. A forward kernel $q(x_t | x_{t-1})$ progressively corrupts x_0 into a noisy sample x_t . A neural network $f_{\theta,t}(x_t)$, parameterized by θ , is trained to approximate the reverse process $p_{\theta}(x_{t-1} | x_t)$. More details are in Appendix B. Certain parameterizations of the diffusion process (Karras et al., 2022; 2024b; Song et al., 2021a) allow for partial or full inversion of a given input sample, producing a noisy sample x_t that preserves the semantic information of x_0 (Zhou et al., 2024a). Hereafter, we denote the inversion process of a sample by F_{inv} and the corresponding decoding process as F_{dec} .

3.2 THE PROCEDURE OF CANONICAL LATENT REPRESENTATION IDENTIFIER (CLARID)

Given a sample x_0 as the seed, which belongs to class condition c , we first invert it to timestep t_e via F_{inv} . Denote x_{t_e} as the latent code of x_0 at t_e . In this latent space, we identify a set of extraneous directions that preserve class identity yet induce large changes. Concretely, let $f_{\theta}(\cdot)$ be the CDM’s feature extractor at layer l and timestep t_e . A first-order Taylor expansion around x_{t_e} gives:

$$f_{\theta,t_e}^l(x_{t_e} + v) \approx f_{\theta,t_e}^l(x_{t_e}) + \nabla f_{\theta,t_e}^l(x_{t_e}) \cdot v = f_{\theta,t_e}^l(x_{t_e}) + J_{\theta,t_e}^l(x_{t_e}) \cdot v. \quad (1)$$

$J_{\theta,t_e}^l(x_{t_e})$ denotes the Jacobian of $f_{\theta,t_e}^l(x_{t_e})$. Hereafter, we drop θ and l to avoid clutter. A vector v that can cause large changes in the output carries some semantic information (Park et al., 2023a; Song et al., 2023), but such information is **not necessarily class-relevant**, as shown in Figure 1. The change caused by vector v is defined as the L2-norm of the Jacobian vector product: $\|J_{t_e}(x_{t_e}) \cdot v\|_2$. Accordingly, the directions that lead to large changes in the output are the right singular vectors of J_{t_e} . When f receives conditional information c , modifying the latent code along those directions tend to **preserve the class identity** while primarily altering the appearance of the decoded sample. By construction, these directions capture high-variance aspects of the image that the CDM can vary freely with the class condition unchanged. We show examples in Appendix D. We therefore remove the components of x_{t_e} along those extraneous directions to obtain CanoRep \tilde{x}_{t_e} , by projecting x_{t_e} onto the subspace which is orthogonal to these directions, as described in Eq.2.

$$\tilde{x}_{t_e} := x_{t_e, \perp V_k} = (\mathbf{I} - \mathbf{V}_k \mathbf{V}_k^T) x_{t_e}. \quad (2)$$

Here, $\mathbf{V} = [v_1, v_2, \dots]$ is the right singular vector matrix of J_{t_e} , and k denotes the number of top singular vectors to be removed. We then apply F_{dec} to \tilde{x}_{t_e} to obtain one Canonical Sample \tilde{x}_0 for the input condition c . Note that the CDM is conditioned on c in the whole procedure. Additional discussion on the choice of the layer index l for computing J_{t_e} is provided in Appendix E.

3.2.1 FINDING APPROPRIATE t_e AND k

Choosing an appropriate t_e is critical for CLARID to be effective. If t_e is too small, e.g., $t_e \approx \frac{T}{100}$, the synthesized sample barely changes and remains close to the original input. On the other hand, when t_e is too large, the conditional signal becomes ineffective and fails to guide the generation (Kynkäänniemi et al., 2024; Zhang et al., 2022). To find the largest time step where the conditional signal is still able

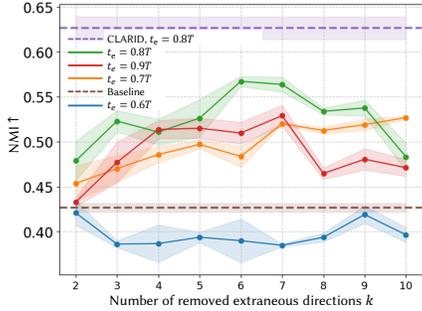


Figure 3: The normalized mutual information (NMI, higher is better) between cluster assignments of CDM features and the ground truth labels. CLARID achieves the highest NMI. Baseline is the original CDM features.

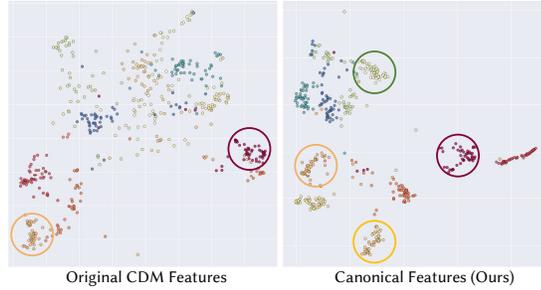


Figure 4: A 2D UMAP (McInnes et al., 2018) projection of the CDM feature space, showing clusters for ten classes. Colors indicate classes. CLARID yields more compact feature clusters than the original samples (green, yellow) and preserve the existing ones (red, orange).

to steer the CDM’s output toward the desired class, we perform a two-stage sampling process starting from $p(\mathbf{x}_T) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$:

1. **Unconditional stage** ($T \leq t < t_e$): Forward the model without class conditioning, *i.e.*, $\mathbf{c} = \emptyset$.
2. **Conditional stage** ($t_e \leq t \leq 0$): The class condition is introduced and retained down to $t = 0$.

We generate m samples for each condition in a given dataset, and measure the accuracy of pre-trained classifiers on the generated samples. We find the saturation point of the accuracy curve as the appropriate t_e . We show the accuracy curve and more details for an ImageNet 256×256 -pretrained DiT (Peebles & Xie, 2023) and a Stable Diffusion model (Rombach et al., 2022) in Appendix H.1.

We fix the inversion time step t_e across all samples, but the number of directions to discard should be adapted to each sample. To decide k , we examine how strongly the top- k extraneous directions alter the sample’s visual appearance, quantified by the explained variance ratio (EVR) $S_k = \sum_{i=1}^k \sigma_i^2 / \sum_{j=1}^n \sigma_j^2$, where σ_i is the i th singular value of \mathbf{J}_{t_e} and n is a hyperparameter. Intuitively, the extraneous direction that leads to larger variations carries less core class semantics. We compute a sequence of the EVR S_1, S_2, \dots, S_n and set the elbow point of the sequence to be the optimal k for a given sample. Compared to a fixed k , the adaptive choice will find the point where the effect of the extraneous directions diminishes. Visual examples illustrating the necessity of adapting k for each sample, along with details of deciding n and the calculation of the elbow point, are provided in Appendix H.2. It is worth noting that all hyperparameters in CLARID are **model-level** instead of sample-level. The hyperparameter selection takes less than 5h for DiT and 10h for SD.

Empirical validation. To demonstrate the effectiveness of our strategies on finding t_e and k , we quantify the feature quality of a CDM and show that our strategies achieve the best one. Intuitively, by retaining the essential class features with minimal class-irrelevant information, the CDM features associated with the CanoReps (Canonical Features) should be more easily separable in the feature space. To show this, we first perform K-means clustering on the CDM features of 1000 samples from 20 different ImageNet classes, as well as their corresponding Canonical Features. We then quantify the feature quality using normalized mutual information (NMI) between the cluster assignments and the ground truth class labels. This method is training-free, hence efficient. We perform the analysis on an ImageNet 256×256 -pretrained DiT-XL model (Peebles & Xie, 2023). **The result is shown in Figure 3: each curve corresponds to a fixed inversion time step $t_e \in \{0.6T, 0.7T, 0.8T, 0.9T\}$, and the x-axis varies the number k of removed extraneous directions; the "Baseline" curve shows the NMI of the original CDM features without removing any direction. The figure also plots error bars, which denote the standard deviation over three independent runs, each using a different set of 20 ImageNet classes.** Removing different numbers of extraneous directions, *i.e.*, varying k , yields different NMI scores, whereas CLARID adaptively chooses k for each sample and achieves the highest NMI. Qualitatively, CLARID produces feature clusters that are more compact than those formed by the original samples while preserving the existing clusters, as shown in Figure 4. The results of Stable Diffusion model in Appendix H.2.1 show a similar trend as DiT, which further validates our conclusion. **We refer the readers to Appendix G and H for the details of the ImageNet20 dataset and experiment setup.**

3.3 A PROOF-OF-CONCEPT EXPERIMENT

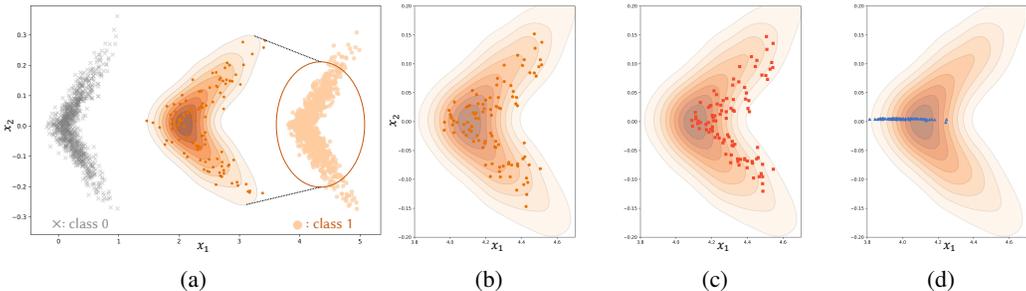


Figure 5: A toy example of CLARID. (a): The samples of class 0 and class 1; (b,c,d): The generated class-1 samples of (b): Plain F_{dec} , (c): CFG, and (d): CLARID, respectively, after applying F_{inv} to the samples. CLARID produces Canonical Samples that lie on a 1D manifold inside class 1, offering an intuitive visual summary of the core class semantics.

We demonstrate the effect of CLARID with a simple yet illustrative example. Figure 5a shows samples generated from our toy generative process and the corresponding density map. The process first generates class-specific samples on a segment $L = \{(x_1, 0) | 4y - 0.1 \leq x_1 \leq 4y + 0.1\}$, where $y \in \{0, 1\}$ denotes the class labels. It then adds class-independent noise to the points to generate the observed data. The samples from class 0 are included solely to introduce an inter-class contrast. We train a small class-conditional diffusion model on this data. After training, we perform CLARID on randomly selected samples from class 1. As a baseline, we take the latent codes obtained with F_{inv} and perform classifier-free guidance (CFG), steering the generation process toward regions with higher class-1 likelihood. The results are shown in Figure 5. Notably, CLARID pushes most samples to a 1D manifold inside class 1, whereas CFG mainly steers the samples away from class 0. This low-dimensional manifold described by Canonical Samples can be regarded as a summarization of class 1 information in this case. The underlying structure revealed by Canonical Samples corresponds to one of the true generative processes for the data, which is the one used here. We refer readers to Section K in the Appendix for details of the toy experiments. Reliably recovering the exact generative model is intractable due to the identifiability issue (Locatello et al., 2019). The solution generally requires extra inductive bias in modeling data distribution, which we leave for future work.

3.4 QUALITATIVE RESULTS

Scaling up, we demonstrate the qualitative effect of CLARID for two CDMs, a class-conditioned DiT (Peebles & Xie, 2023) trained on the ImageNet 256×256 dataset and a Stable Diffusion 2.1 model (Rombach et al., 2022) trained on a subset of LAION-5B (Schuhmann et al., 2022). We use DDIM (Song et al., 2021a) as the sampler and use 100 diffusion steps for both inversion and decoding. Visual results are shown in Figure 6 (more in Appendix Q). We adopt classifier-free guidance (CFG) after removing extraneous directions in CLARID to ensure the data fidelity, and compare the results against pure CFG. The CanoReps visualized as Canonical Samples show that our method preserves the core class information in the original images, while CFG focuses on increasing the class-conditional likelihood. Occasionally, CLARID can select suboptimal t_e and k , leading to artefacts in the generated images. The failure cases are discussed in Appendix H.3.

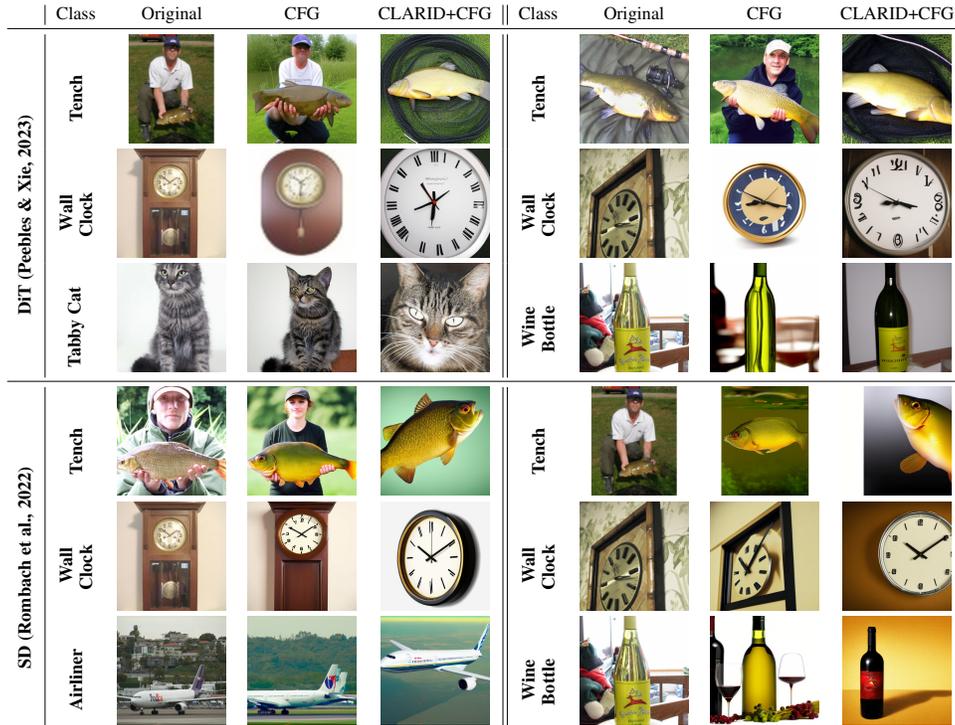
3.5 ON THE GENERALIZATION OF CLARID

While we focus on class-conditional DMs to develop the CLARID framework, it also extends naturally to text-conditioned models e.g. Stable Diffusion (Rombach et al., 2022), as shown in Figure 6 and Appendix F. Text prompts span a far richer semantic space than one-hot class labels, giving finer control over where the CanoReps lie. We show examples of the effect using different prompts on the same input in Appendix F.1. Understanding how this semantic structure shapes the located CanoReps is an interesting direction to explore. CLARID is also compatible with different DM samplers, as shown in Appendix F.2.

4 THE APPLICATION OF CANOREPS

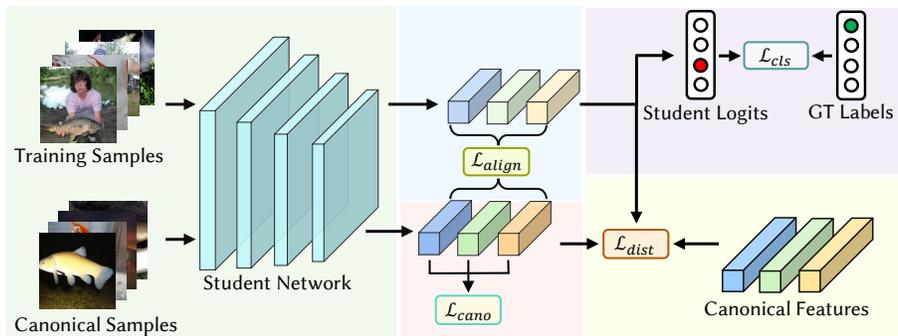
As demonstrated in Section 3, CanoReps correspond to the essential class information learned by the CDMs. Building on this insight, we design a feature distillation framework for CanoReps, termed

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349



350 Figure 6: Comparison of classifier-free guidance (CFG) and CLARID on an ImageNet 256×256 DiT
351 (Peebles & Xie, 2023) and Stable Diffusion 2.1 (SD) (Rombach et al., 2022). We use the following
352 prompt template for SD: *a photo of <Class>* (Li et al., 2023a). CLARID focuses on identifying
353 CanoReps to preserve the core class information, yielding Canonical Samples that provide an
354 interpretable summary of the essential class semantics, whereas CFG aims at finding high-likelihood
355 images. Two Canonical Samples from the *Tench* and *Wall Clock* class are presented to show that
356 CanoReps do not collapse to a single constant vector. All "Original" images are taken from ImageNet.
357 More visual results are in Appendix Q.

358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377



373 Figure 7: Overview of *CaDistill*. We align the student's features of training images to those of
374 Canonical Samples using \mathcal{L}_{align} . The student is trained to discriminate between Canonical Samples
375 in different classes by optimizing \mathcal{L}_{cano} . The CDM, as a teacher, provides Canonical Features for
376 feature distillation with \mathcal{L}_{dist} . Finally, the student is supervised on the ground-truth labels via \mathcal{L}_{cls} .

CaDistill, as illustrated in Figure 7. Given a training batch $B = \{(\mathbf{x}_1, \mathbf{c}_1), (\mathbf{x}_2, \mathbf{c}_2), \dots, (\mathbf{x}_b, \mathbf{c}_b)\}$ where \mathbf{x}_i is an image and \mathbf{c}_i is the class label, we first **randomly** select a CanoRep $\tilde{\mathbf{x}}_i$ corresponding to the category of each sample. The random sampling relaxes the constraint of the size equivalence between the training dataset and the set of CanoReps. Let $P_i = \{j \in \{1, 2, \dots, b\} \mid \mathbf{c}_j = \mathbf{c}_i\}$ be the set of indices of samples in the batch with class label \mathbf{c}_i . We then compute student features $\mathbf{z}_i = g_\phi(\mathbf{x}_i) \in \mathbb{R}^d$, $\tilde{\mathbf{z}}_i = g_\phi(\tilde{\mathbf{x}}_i) \in \mathbb{R}^d$, where g denotes the student network and ϕ is the set of its training parameters. The first component in **CaDistill**, \mathcal{L}_{align} , is designed to pull each training image feature towards all same-class Canonical Samples, and push away the ones from other classes. It is given in Eq. 3.

$$\mathcal{L}_{align} = -\frac{1}{b} \sum_{i=1}^b \frac{1}{|P_i|} \sum_{j \in P_i} \log \frac{\exp(\mathbf{z}_i \cdot \tilde{\mathbf{z}}_j / \tau)}{\sum_{k=1}^B \exp(\mathbf{z}_i \cdot \tilde{\mathbf{z}}_k / \tau)}. \quad (3)$$

τ is a temperature hyperparameter. Second, we encourage all Canonical Samples of the same class to cluster in the student’s feature space, and separate those of different classes. It is achieved by minimizing \mathcal{L}_{cano} in Eq. 4. When a Canonical Sample has no same-class positives, we simply optimize the denominator. See Appendix M.7.5 for more details.

$$\mathcal{L}_{cano} = -\frac{1}{b} \sum_{i=1}^b \frac{1}{|P_i| - 1} \sum_{j \in P_i, j \neq i} \log \frac{\exp(\tilde{\mathbf{z}}_i \cdot \tilde{\mathbf{z}}_j / \tau)}{\sum_{k \neq i} \exp(\tilde{\mathbf{z}}_i \cdot \tilde{\mathbf{z}}_k / \tau)}. \quad (4)$$

We perform an ablation study on the design of \mathcal{L}_{cano} by replacing it with a classification loss on the Canonical Samples in Appendix M.7.5, demonstrating the advantage of the current \mathcal{L}_{cano} . For feature distillation, we transfer the structures of Canonical Features to the student via minimizing \mathcal{L}_{dist} . Specifically, we use the Centered Kernel Alignment (CKA) (Kornblith et al., 2019) metric to align the student’s representations of both the training images and the Canonical Samples with the Canonical Features extracted from the CDM. Maximizing CKA, *i.e.* minimizing \mathcal{L}_{dist} , aligns the linear subspace spanned by the student’s feature vectors with that of the teacher, effectively distilling the teacher’s class-discriminative structure into the student. Denote the student feature matrices of training images and Canonical Samples as $\mathbf{Z} \in \mathbb{R}^{b \times d}$ and $\tilde{\mathbf{Z}} \in \mathbb{R}^{b \times d}$, respectively, and the Canonical Feature matrix as $\mathcal{A} \in \mathbb{R}^{b \times d'}$. The Canonical Features are extracted from a frozen CDM for all $\tilde{\mathbf{x}}_i$ in a training batch.

$$\mathcal{L}_{dist} = \lambda_{cka} \log(1 - \text{CKA}(\mathbf{Z}, \mathcal{A})) + (1 - \lambda_{cka}) \log(1 - \text{CKA}(\tilde{\mathbf{Z}}, \mathcal{A})). \quad (5)$$

We find that using CKA for structural alignment outperforms existing diffusion-based feature distillation methods (Li et al., 2023b; Park et al., 2019; Romero et al., 2015; Yang & Wang, 2023) and refer the readers to Appendix M.1 for more details. Finally, the student is trained with the standard cross-entropy loss \mathcal{L}_{cls} on ground-truth labels. The final loss function for **CaDistill** is given in Eq. 6.

$$\mathcal{L}_{CaDistill} = \mathcal{L}_{cls} + \lambda_{cs}(\lambda_{cf} \mathcal{L}_{align} + (1 - \lambda_{cf}) \mathcal{L}_{cano}) + \lambda_{dist} \mathcal{L}_{dist} \quad (6)$$

4.1 CANOREPS IN PRACTICE

We conduct experiments of **CaDistill** on CIFAR10 (Krizhevsky et al., 2009) using a pre-trained UNet-based CDM (Ho et al., 2020; Xiang et al., 2023) and on ImageNet using the ImageNet 256 × 256-trained Diffusion Transformer (DiT) (Peebles & Xie, 2023). We choose two different DM architectures to demonstrate the applicability of our method. The CDM on CIFAR10 does not have an unconditional branch. The DMs are trained on the same dataset as the student models, hence no extra data is considered (Shama et al., 2024). We follow this principle to choose the baselines, comparing our method with the SOTA diffusion-based feature distillation (Yang & Wang, 2023) and data augmentation (Shama et al., 2024) methods. In addition, we design two important baselines. The first one, DMDistill, is to distill the raw space structure in CDMs using \mathcal{L}_{dist} on all training samples, which represents the current mainstream idea of using DMs as teachers in feature distillation. It outperforms existing diffusion-based feature distillation losses (Li et al., 2023b; Park et al., 2019; Romero et al., 2015; Yang & Wang, 2023) (Appendix M.1). Second, for the CFGDistill baseline, we train the student model using the same framework of **CaDistill**, except that the images and features are obtained using CFG. For the CDM on CIFAR10, we sample new images for CFGDistill as this CDM lacks the unconditional branch and cannot perform CFG. We fix t_r , *i.e.* the feature extraction time step, for all methods that do not adaptively change it (Yang & Wang, 2023). For the student network, we use ResNet18 (He et al., 2016) on CIFAR10 and ResNet152 (He et al., 2016) on ImageNet, two well-established convolutional neural network baselines (Yang & Wang, 2023). ResNet50 results on

Table 1: Quantitative comparisons of *CaDistill* and baselines on CIFAR-10 (Krizhevsky et al., 2009) (ResNet-18) and ImageNet (Deng et al., 2009) (ResNet-152). \mathcal{D} : Dataset. Adversarial robustness benchmarks: PGD (Madry et al., 2018), CW (Carlini & Wagner, 2017), APGD-DLR / APGD-CE (Croce & Hein, 2020); Evaluations of generalization ability : Corruption (CIFAR10-C and ImageNet-C) (Hendrycks & Dietterich, 2018), ImageNet-A (Djolonga et al., 2021), ImageNet-Real (Beyer et al., 2020), CIFAR10.1 (Recht et al., 2018), CIFAR10.2 (Lu et al., 2020). Data_{DM} is the portion of data for which the DM acts as teacher. Higher is better. Values lower than the vanilla model are in red. †: the model relies on unconditional DMs, or the training cannot be performed, see Appendix M.6.

\mathcal{D}	Model	Data_{DM}	Clean	PGD	CW	APGD-DLR	APGD-CE	Corruption	C10.1/IM-A	C10.2/IM-Real
CIFAR10	Vanilla	/	92.4	33.4	20.9	34.2	32.0	76.1	82.3	78.0
	SupCon (2020)	/	92.7	29.1	16.8	34.8	29.9	76.9	81.8	78.3
	RepFusion† (2023)	100%	92.7	30.3	17.2	32.1	29.2	75.3	83.4	78.5
	DMDistill	100%	92.9	41.3	32.8	38.7	36.0	76.7	83.2	79.1
	CFGDistill	10%	92.9	43.7	37.3	40.9	39.0	76.6	83.8	79.5
	<i>CaDistill</i>	10%	93.1	47.9	43.1	44.1	43.3	77.7	84.5	79.7
ImageNet	Vanilla	/	79.3	22.0	20.7	24.7	23.3	54.2	14.1	85.5
	DiffAug (2024)	100%	79.6	24.5	21.7	26.3	25.7	55.5	13.0	85.6
	DMDistill	100%	79.4	23.8	23.1	25.7	24.9	51.8	12.5	85.6
	CFGDistill	10%	79.1	27.6	30.8	30.2	28.2	54.1	13.7	85.6
	<i>CaDistill</i>	10%	79.5	29.7	32.2	32.5	29.6	55.1	14.9	85.8

ImageNet are in Appendix M, showing a similar trend as ResNet152. We evaluate the adversarial robustness as well as the generalization ability, including in-distribution and out-of-distribution, of the student. We refer readers to Appendix L for more training and evaluation procedure details. For all *CaDistill* experiments, CLARID is used as a single, offline preprocessing step. This avoids any online computational burden, making it efficient for scaling up. More details of the computational costs and scalability of CLARID are in Appendix I.

The results are shown in Table 1. *CaDistill* consistently improves the adversarial robustness and generalization ability of the student, while the main effect of SOTA methods is improving the clean accuracy or a single aspect of the robustness. A more detailed discussion of the results, including the significance of such a consistent improvement and the importance of both our Canonical Samples and loss functions, is in Appendix N. Moreover, we consider an interesting and challenging baseline that is capable of extracting class semantics, *i.e.*, using the class token in an ImageNet-pretrained ViT (Touvron et al., 2021; 2022) as the teacher. We show the results in Appendix M.4. *CaDistill* is effective when the student has transformer-based vision backbones, as shown in Appendix M.8.

The student, trained with *CaDistill*, focuses more on the core class signal. We demonstrate this by a test on the Backgrounds Challenge (Xiao et al.), where the background of an image that is irrelevant to the class identity is either removed or shuffled. The results are given in Table 2. *CaDistill* improves the student performance on BG-Rand and Only-FG while maintaining the accuracy on the Original and BG-Same splits, indicating a mitigation in the model’s dependence on the spurious background cues for classification. A more detailed discussion of the results is in Appendix M.5.

Table 2: Results on the Backgrounds Challenge (Xiao et al.). Higher is better. See Appendix M.5 for details.

Model	Original	BG-Same	BG-Rand	Only-FG
Vanilla	96.9	91.3	85.6	89.6
DiffAug (2024)	97.0	90.5	85.2	89.2
DMDistill	97.0	90.0	84.6	88.1
CFGDistill	97.0	91.5	85.4	89.4
<i>CaDistill</i>	97.0	91.5	86.5	90.4

Ablation studies. Appendix M.7 include the ablation analysis on *CaDistill*, including:

- **Number of CanoReps.** We demonstrate that 10% of CanoReps is sufficient for achieving competitive performance, implying the low-dimensionality property of the class manifolds in CDMs.
- **Necessity of \mathcal{L}_{align} , \mathcal{L}_{cano} , \mathcal{L}_{dist} , and their balance.** We empirically conclude the effects of all loss functions and choose the optimal weighting schemes, including λ_{cs} , λ_{cf} , λ_{dist} , λ_{cka} .
- **CFG magnitude.** We perform CFG after obtaining CanoReps, and show a proper choice of its magnitude. Importantly, a higher CFG magnitude does not contribute to better performance,

486 indicating that *CaDistill* is not merely providing a converging prior on the student features. We
487 also provide a discussion on this topic in Appendix G.3.
488

489 5 DISCUSSION AND LIMITATION 490

491 **Discussion.** Our work offers a fresh perspective on representation learning: rather than designing
492 objectives to enforce invariance or discriminative properties, we *subtract* non-discriminative com-
493 ponents from DM’s rich feature space to reveal CanoReps that encode core class semantics. This
494 complements the current representation learning research. We further identify several properties,
495 including its relationship to generative similarity (Marjeh et al., 2024), and potential applications of
496 CLARID in Appendix O.
497

498 **Limitation.** CLARID has certain limitations. It can occasionally select a suboptimal projection time
499 step, t_e , or the total number of extraneous directions considered, n , as discussed in Section 3.4 and
500 Appendix Section H. We discuss potential solutions to those issues in Appendix H.3. Calculating the
501 singular vectors of the Jacobian of CDMs is computationally intensive. Whether the application of
502 CanoReps can be effective on larger-scale problems, *e.g.*, ImageNet22K, remains a question.
503

504 6 CONCLUSION 505

506 We introduce Canonical LATent Representation IDENTIFIER (CLARID), a principled method to uncover
507 the core categorical information encoded in pre-trained conditional diffusion models (CDMs). By
508 removing extraneous directions from a sample’s latent code, CLARID produces *CanoRep*, whose
509 internal feature—Canonical Feature—distills the class-defining semantics of each category. Decoding
510 CanoReps yields Canonical Samples that offer an interpretable and compact summary of the class.
511 Quantitatively, Canonical Features form more compact and easily separable clusters in CDM feature
512 space than the original inputs. Building on CanoReps, we have proposed *CaDistill*, a diffusion-based
513 feature distillation framework. The teacher CDM transfers core class semantics to the student only via
514 CanoReps, which is equivalent to merely 10% of the original data, while the student is being trained on
515 the full training set. The student achieves strong adversarial robustness and generalization, focusing
516 more on the true class signal instead of spurious background cues than the original model. Together,
517 CLARID and *CaDistill* demonstrate that CDMs can be transformed from black-box generators into
518 compact, interpretable teachers for robust representation learning.
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

7 REPRODUCIBILITY STATEMENT

We disclose all the method and experiment details regarding CLARID in Section 3.2, Appendix E, G, H. For giving an intuitive understanding of CLARID, we provide a 2D toy model whose details are in Appendix K. The application of CLARID is *CaDistill*, for which we provide detailed training and evaluation settings in Appendix L.

8 ETHICS STATEMENT

The image generators used in our experiments may sometimes contain wrong information. That said, we are not focusing on generating high-fidelity images or improving the synthesizing quality.

REFERENCES

- Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square attack: a query-efficient black-box adversarial attack via random search. In *European Conference on Computer Vision*, pp. 484–501. Springer, 2020.
- Shekoofeh Azizi, Simon Kornblith, Chitwan Saharia, Mohammad Norouzi, and David J Fleet. Synthetic data from diffusion models improves imagenet classification. *Transactions on Machine Learning Research*.
- Hritik Bansal and Aditya Grover. Leaving reality to imagination: Robust classification via generated datasets. In *ICLR 2023 Workshop on Trustworthy and Reliable Large-Scale Machine Learning Models*.
- Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth words: A vit backbone for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22669–22679, 2023.
- Dmitry Baranchuk, Andrey Voynov, Ivan Rubachev, Valentin Khulkov, and Artem Babenko. Label-efficient semantic segmentation with diffusion models. In *International Conference on Learning Representations*, 2022.
- Lucas Beyer, Olivier J Hénaff, Alexander Kolesnikov, Xiaohua Zhai, and Aäron van den Oord. Are we done with imagenet? *arXiv preprint arXiv:2006.07159*, 2020.
- Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 39–57. IEEE, 2017.
- Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. *ACM Transactions on Graphics (TOG)*, 42(4):1–10, 2023.
- Ricky TQ Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disentanglement in variational autoencoders. *Advances in Neural Information Processing Systems*, 31, 2018.
- Siyi Chen, Huijie Zhang, Minzhe Guo, Yifu Lu, Peng Wang, and Qing Qu. Exploring low-dimensional subspace in diffusion models for controllable image editing. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024a.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning*, pp. 1597–1607. PmlR, 2020a.
- Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey E Hinton. Big self-supervised models are strong semi-supervised learners. *Advances in Neural Information Processing Systems*, 33:22243–22255, 2020b.
- Xinlei Chen, Zhuang Liu, Saining Xie, and Kaiming He. Deconstructing denoising diffusion models for self-supervised learning. *arXiv preprint arXiv:2401.14404*, 2024b.

- 594 Yixin Cheng, Grigorios Chrysos, Markos Georgopoulos, and Volkan Cevher. Multilinear operator
595 networks. In *The Twelfth International Conference on Learning Representations*, 2023.
- 596
- 597 Edo Collins, Radhakrishna Achanta, and Sabine Susstrunk. Deep feature factorization for concept
598 discovery. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 336–352,
599 2018.
- 600 Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of
601 diverse parameter-free attacks. In *International Conference on Machine Learning*, pp. 2206–2216.
602 PMLR, 2020.
- 603 Joel Dapello, Kohitij Kar, Martin Schrimpf, Robert Baldwin Geary, Michael Ferguson, David Daniel
604 Cox, and James J DiCarlo. Aligning model and macaque inferior temporal cortex representations
605 improves model-to-human behavioral alignment and adversarial robustness. In *The Eleventh
606 International Conference on Learning Representations*, 2023.
- 607
- 608 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-
609 scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern
610 Recognition*, pp. 248–255. IEEE, 2009.
- 611 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances
612 in Neural Information Processing Systems*, 34:8780–8794, 2021.
- 613
- 614 Josip Djolonga, Jessica Yung, Michael Tschannen, Rob Romijnders, Lucas Beyer, Alexander
615 Kolesnikov, Joan Puigcerver, Matthias Minderer, Alexander D’Amour, Dan Moldovan, et al. On
616 robustness and transferability of convolutional neural networks. In *Proceedings of the IEEE/CVF
617 Conference on Computer Vision and Pattern Recognition*, pp. 16458–16468, 2021.
- 618 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
619 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image
620 is worth 16x16 words: Transformers for image recognition at scale. In *International Conference
621 on Learning Representations*, 2020.
- 622 Arthur Douillard, Alexandre Ramé, Guillaume Couairon, and Matthieu Cord. DyTox: Transformers
623 for continual learning with dynamic token expansion. In *Proceedings of the IEEE/CVF Conference
624 on Computer Vision and Pattern Recognition*, pp. 9285–9295, 2022.
- 625
- 626 Yunxiang Fu, Chaoqi Chen, Yu Qiao, and Yizhou Yu. Dreamda: Generative data augmentation with
627 diffusion models. *arXiv preprint arXiv:2403.12803*, 2024.
- 628 Paul Gavrikov and Janis Keuper. Can biases in imagenet models explain generalization? In
629 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
630 22184–22194, 2024.
- 631
- 632 Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and
633 Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves
634 accuracy and robustness. In *International Conference on Learning Representations*, 2018a.
- 635
- 636 Robert Geirhos, Carlos RM Temme, Jonas Rauber, Heiko H Schütt, Matthias Bethge, and Felix A
637 Wichmann. Generalisation in humans and deep neural networks. *Advances in Neural Information
638 Processing Systems*, 31, 2018b.
- 639
- 640 Sven Gowal, Sylvestre-Alvise Rebuffi, Olivia Wiles, Florian Stimberg, Dan Andrei Calian, and
641 Timothy A Mann. Improving robustness using generated data. *Advances in Neural Information
642 Processing Systems*, 34:4218–4233, 2021.
- 643
- 644 René Haas, Inbar Huberman-Spiegelglas, Rotem Mulayoff, Stella Graßhof, Sami S Brandt, and
645 Tomer Michaeli. Discovering interpretable directions in the semantic latent space of diffusion
646 models. In *2024 IEEE 18th International Conference on Automatic Face and Gesture Recognition
647 (FG)*, pp. 1–9. IEEE, 2024.
- 648
- 649 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
650 recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,
651 pp. 770–778, 2016.

- 648 Ruifei He, Shuyang Sun, Xin Yu, Chuhui Xue, Wenqing Zhang, Philip Torr, Song Bai, and Xiaojuan
649 Qi. Is synthetic data from generative models ready for image recognition? In *The Eleventh
650 International Conference on Learning Representations*, 2023.
- 651 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common
652 corruptions and perturbations. In *International Conference on Learning Representations*, 2018.
- 653 Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or. Prompt-
654 to-prompt image editing with cross-attention control. In *The Eleventh International Conference on
655 Learning Representations*.
- 656 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on
657 Deep Generative Models and Downstream Applications*.
- 658 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in
659 Neural Information Processing Systems*, 33:6840–6851, 2020.
- 660 Taicheng Huang, Zonglei Zhen, and Jia Liu. Semantic relatedness emerges in deep convolutional
661 neural networks designed for object recognition. *Frontiers in Computational Neuroscience*, 15:
662 625804, 2021.
- 663 Khawar Islam, Muhammad Zaigham Zaheer, Arif Mahmood, and Karthik Nandakumar. Diffusemix:
664 Label-preserving data augmentation with diffusion models. In *Proceedings of the IEEE/CVF
665 Conference on Computer Vision and Pattern Recognition*, pp. 27621–27630, 2024.
- 666 Priyank Jaini, Kevin Clark, and Robert Geirhos. Intriguing properties of generative classifiers. In
667 *The Twelfth International Conference on Learning Representations*.
- 668 Jaeseok Jeong, Mingi Kwon, and Youngjung Uh. Training-free content injection using h-space
669 in diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of
670 Computer Vision*, pp. 5151–5161, 2024.
- 671 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-
672 based generative models. *Advances in Neural Information Processing Systems*, 35:26565–26577,
673 2022.
- 674 Tero Karras, Miika Aittala, Tuomas Kynkäänniemi, Jaakko Lehtinen, Timo Aila, and Samuli Laine.
675 Guiding a diffusion model with a bad version of itself. *Advances in Neural Information Processing
676 Systems*, 37:52996–53021, 2024a.
- 677 Tero Karras, Miika Aittala, Jaakko Lehtinen, Janne Hellsten, Timo Aila, and Samuli Laine. Analyzing
678 and improving the training dynamics of diffusion models. In *Proceedings of the IEEE/CVF
679 Conference on Computer Vision and Pattern Recognition*, pp. 24174–24184, 2024b.
- 680 Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron
681 Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in Neural
682 Information Processing Systems*, 33:18661–18673, 2020.
- 683 Hoki Kim. Torchattacks: A pytorch repository for adversarial attacks. *arXiv preprint
684 arXiv:2010.01950*, 2020.
- 685 Yeongmin Kim, Kwanghyeon Lee, Minsang Park, Byeonghu Na, and Il-chul Moon. Diffusion bridge
686 autoencoders for unsupervised representation learning. In *The Thirteenth International Conference
687 on Learning Representations*.
- 688 Yunji Kim, Jiyoung Lee, Jin-Hwa Kim, Jung-Woo Ha, and Jun-Yan Zhu. Dense text-to-image
689 generation with attention modulation. In *Proceedings of the IEEE/CVF International Conference
690 on Computer Vision*, pp. 7701–7711, 2023.
- 691 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. 2015.
- 692 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint
693 arXiv:1312.6114*, 2014.

- 702 Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural
703 network representations revisited. In *International Conference on Machine Learning*, pp. 3519–
704 3529. PMLR, 2019.
- 705 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 706
707 Mingi Kwon, Jaeseok Jeong, and Youngjung Uh. Diffusion models already have a semantic latent
708 space. In *International Conference on Learning Representations*, 2023.
- 709
710 Tuomas Kynkäänniemi, Miika Aittala, Tero Karras, Samuli Laine, Timo Aila, and Jaakko Lehtinen.
711 Applying guidance in a limited interval improves sample and distribution quality in diffusion
712 models. In *Conference on Neural Information Processing Systems*, 2024.
- 713
714 Alexander C Li, Mihir Prabhudesai, Shivam Duggal, Ellis Brown, and Deepak Pathak. Your diffusion
715 model is secretly a zero-shot classifier. In *Proceedings of the IEEE/CVF International Conference*
716 *on Computer Vision*, pp. 2206–2217, 2023a.
- 717
718 Daiqing Li, Huan Ling, Amlan Kar, David Acuna, Seung Wook Kim, Karsten Kreis, Antonio Torralba,
719 and Sanja Fidler. Dreamteacher: Pretraining image backbones with deep generative models. In
720 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 16698–16708,
721 2023b.
- 722
723 Bingyan Liu, Chengyu Wang, Tingfeng Cao, Kui Jia, and Jun Huang. Towards understanding cross
724 and self-attention in stable diffusion for text-guided image editing. In *Proceedings of the IEEE/CVF*
Conference on Computer Vision and Pattern Recognition, pp. 7817–7826, 2024.
- 725
726 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo.
727 Swin Transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the*
728 *IEEE/CVF International Conference on Computer Vision*, pp. 10012–10022, 2021.
- 729
730 Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng
731 Zhang, Li Dong, et al. Swin Transformer v2: Scaling up capacity and resolution. In *Proceedings of*
732 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12009–12019, 2022.
- 733
734 Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard Schölkopf,
735 and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled
736 representations. In *International Conference on Machine Learning*, pp. 4114–4124. PMLR,
737 2019.
- 738
739 Shangyun Lu, Bradley Nott, Aaron Olson, Alberto Todeschini, Hossein Vahabi, Yair Carmon, and
740 Ludwig Schmidt. Harder or different? a closer look at distribution shift in dataset reproduction. In
741 *ICML Workshop on Uncertainty and Robustness in Deep Learning*, volume 5, pp. 15, 2020.
- 742
743 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
744 Towards deep learning models resistant to adversarial attacks. In *International Conference on*
745 *Learning Representations*, 2018.
- 746
747 Raja Marjeh, Sreejan Kumar, Declan Campbell, Liyi Zhang, Gianluca Bencomo, Jake Snell, and
748 Thomas L Griffiths. Learning human-aligned representations with contrastive learning and genera-
749 tive similarity. *arXiv preprint arXiv:2405.19420*, 2024.
- 750
751 Leland McInnes, John Healy, Nathaniel Saul, and Lukas Großberger. Umap: Uniform manifold
752 approximation and projection. *Journal of Open Source Software*, 3(29):861, 2018.
- 753
754 Benyuan Meng, Qianqian Xu, Zitai Wang, Xiaochun Cao, and Qingming Huang. Not all diffusion
755 model activations have been evaluated as discriminative features. *Advances in Neural Information*
Processing Systems, 37:55141–55177, 2024.
- 756
757 Mazda Moayeri, Kiarash Banhashem, and Soheil Feizi. Explicit tradeoffs between adversarial
758 and natural distributional robustness. *Advances in Neural Information Processing Systems*, 35:
759 38761–38774, 2022.

- 756 Soumik Mukhopadhyay, Matthew Gwilliam, Vatsal Agarwal, Namitha Padmanabhan, Archana
757 Swaminathan, Srinidhi Hegde, Tianyi Zhou, and Abhinav Shrivastava. Diffusion models beat gans
758 on image classification. *arXiv preprint arXiv:2307.08702*, 2023.
- 759
760 Soumik Mukhopadhyay, Matthew Gwilliam, Yosuke Yamaguchi, Vatsal Agarwal, Namitha Padman-
761 abhan, Archana Swaminathan, Tianyi Zhou, Jun Ohya, and Abhinav Shrivastava. Do text-free
762 diffusion models learn discriminative visual representations? In *European Conference on Computer*
763 *Vision*, pp. 253–272. Springer, 2024.
- 764
765 Tianyu Pang, Kun Xu, Yinpeng Dong, Chao Du, Ning Chen, and Jun Zhu. Rethinking softmax cross-
766 entropy loss for adversarial robustness. In *International Conference on Learning Representations*.
- 767
768 Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In *Proceed-*
769 *ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3967–3976,
770 2019.
- 771
772 Yong-Hyun Park, Mingi Kwon, Jaewoong Choi, Junghyo Jo, and Youngjung Uh. Understanding the
773 latent space of diffusion models through the lens of riemannian geometry. *Advances in Neural*
774 *Information Processing Systems*, 36:24129–24142, 2023a.
- 775
776 Yong-Hyun Park, Mingi Kwon, Junghyo Jo, and Youngjung Uh. Unsupervised discovery of semantic
777 latent directions in diffusion models. *arXiv preprint arXiv:2302.12469*, 2023b.
- 778
779 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*
780 *the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.
- 781
782 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
783 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image
784 synthesis. In *International Conference on Learning Representations*, 2024.
- 785
786 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
787 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
788 models from natural language supervision. In *International Conference on Machine Learning*, pp.
789 8748–8763. PmLR, 2021.
- 790
791 Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishal Shankar. Do cifar-10 classifiers
792 generalize to cifar-10? *arXiv preprint arXiv:1806.00451*, 2018.
- 793
794 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
795 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Confer-*
796 *ence on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.
- 797
798 Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and
799 Yoshua Bengio. Fitnets: Hints for thin deep nets. 2015.
- 800
801 Evgenia Rusak, Lukas Schott, Roland S Zimmermann, Julian Bitterwolf, Oliver Bringmann, Matthias
802 Bethge, and Wieland Brendel. A simple way to make neural networks robust against diverse image
803 corruptions. In *European Conference on Computer Vision*, pp. 53–69. Springer, 2020.
- 804
805 Aninda Saha, Alina Bialkowski, and Sara Khalifa. Distilling representational similarity using centered
806 kernel alignment (cka). In *Proceedings of the the 33rd British Machine Vision Conference (BMVC*
807 *2022)*. British Machine Vision Association, 2022.
- 808
809 Mert Bülent Saryıldız, Karteek Alahari, Diane Larlus, and Yannis Kalantidis. Fake it till you make
810 it: Learning transferable representations from synthetic imagenet clones. In *Proceedings of the*
811 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8011–8021, 2023.
- 812
813 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi
814 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An
815 open large-scale dataset for training next generation image-text models. *Advances in Neural*
816 *Information Processing Systems*, 35:25278–25294, 2022.

- 810 Vikash Sehwal, Saeed Mahlouljifar, Tinashe Handina, Sihui Dai, Chong Xiang, Mung Chiang,
811 and Prateek Mittal. Robust learning meets generative models: Can proxy distributions improve
812 adversarial robustness? In *International Conference on Learning Representations*, 2022.
813
- 814 Sastry Chandramouli Shama, Sri Harsha Dumpala, and Sageev Oore. Diffaug: A diffuse-and-denoise
815 augmentation for training robust classifiers. *Advances in Neural Information Processing Systems*,
816 37:20745–20785, 2024.
- 817 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International
818 Conference on Learning Representations*, 2021a.
819
- 820 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
821 Poole. Score-based generative modeling through stochastic differential equations. In *International
822 Conference on Learning Representations*, 2021b.
- 823 Yue Song, Andy Keller, Nicu Sebe, and Max Welling. Latent traversals in generative models as
824 potential flows. In *Proceedings of the 40th International Conference on Machine Learning*, pp.
825 32288–32303, 2023.
826
- 827 Ajay Subramanian, Elena Sizikova, Najib Majaj, and Denis Pelli. Spatial-frequency channels, shape
828 bias, and adversarial robustness. *Advances in Neural Information Processing Systems*, 36, 2024.
829
- 830 Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé
831 Jégou. Training data-efficient image transformers & distillation through attention. In *International
832 Conference on Machine Learning*, pp. 10347–10357. PMLR, 2021.
- 833 Hugo Touvron, Matthieu Cord, and Hervé Jégou. Deit iii: Revenge of the vit. In *European Conference
834 on Computer Vision*, pp. 516–533. Springer, 2022.
835
- 836 Brandon Trabucco, Kyle Doherty, Max A Gurinas, and Ruslan Salakhutdinov. Effective data
837 augmentation with diffusion models. In *International Conference on Learning Representations*,
838 2024.
- 839 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
840 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing
841 Systems*, 30, 2017.
842
- 843 Peng Wang, Huijie Zhang, Zekai Zhang, Siyi Chen, Yi Ma, and Qing Qu. Diffusion model learns
844 low-dimensional distributions via subspace clustering. In *NeurIPS 2024 Workshop on Mathematics
845 of Modern Machine Learning*.
- 846 Zekai Wang, Tianyu Pang, Chao Du, Min Lin, Weiwei Liu, and Shuicheng Yan. Better diffusion
847 models further improve adversarial training. In *International Conference on Machine Learning*,
848 pp. 36246–36263. PMLR, 2023.
849
- 850 Ross Wightman. Pytorch image models. [https://github.com/rwightman/
851 pytorch-image-models](https://github.com/rwightman/pytorch-image-models), 2019.
- 852 Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement
853 learning. *Machine Learning*, 8:229–256, 1992.
854
- 855 Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei Chen, Zhuang Liu, In So Kweon, and
856 Saining Xie. Convnext v2: Co-designing and scaling convnets with masked autoencoders. In
857 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
858 16133–16142, 2023.
- 859 Zhibiao Wu and Martha Palmer. Verbs semantics and lexical selection. In *Proceedings of the 32nd
860 annual meeting on Association for Computational Linguistics*, pp. 133–138, 1994.
861
- 862 Weilai Xiang, Hongyu Yang, Di Huang, and Yunhong Wang. Denoising diffusion autoencoders are
863 unified self-supervised learners. In *Proceedings of the IEEE/CVF International Conference on
Computer Vision*, pp. 15802–15812, 2023.

- 864 Kai Yuanqing Xiao, Logan Engstrom, Andrew Ilyas, and Aleksander Madry. Noise or signal:
865 The role of image backgrounds in object recognition. In *International Conference on Learning*
866 *Representations*.
- 867 Jiarui Xu, Sifei Liu, Arash Vahdat, Wonmin Byeon, Xiaolong Wang, and Shalini De Mello. Open-
868 vocabulary panoptic segmentation with text-to-image diffusion models. In *Proceedings of the*
869 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2955–2966, 2023.
- 870 Yitao Xu, Tong Zhang, and Sabine Süsstrunk. Adanca: Neural cellular automata as adaptors for
871 more robust vision transformer. In *The Thirty-Eighth Annual Conference on Neural Information*
872 *Processing Systems*, 2024.
- 873 Shipeng Yan, Jiangwei Xie, and Xuming He. DER: Dynamically expandable representation for
874 class incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
875 *Pattern Recognition*, pp. 3014–3023, 2021.
- 876 Xingyi Yang and Xinchao Wang. Diffusion model as representation learner. In *Proceedings of the*
877 *IEEE/CVF International Conference on Computer Vision*, pp. 18938–18949, 2023.
- 878 Zhendong Yang, Zhe Li, Ailing Zeng, Zexian Li, Chun Yuan, and Yu Li. Vitkd: Feature-based
879 knowledge distillation for vision transformers. In *Proceedings of the IEEE/CVF Conference on*
880 *Computer Vision and Pattern Recognition*, pp. 1379–1388, 2024.
- 881 Xufeng Yao, Yuechen ZHANG, Zuyao Chen, Jiaya Jia, and Bei Yu. Distill vision transformers to
882 cnns via low-rank representation approximation. 2022.
- 883 Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu.
884 Coca: Contrastive captioners are image-text foundation models. *Transactions on Machine Learning*
885 *Research*.
- 886 Tan Yu, Xu Li, Yunfeng Cai, Mingming Sun, and Ping Li. S2-mlp: Spatial-shift mlp architecture for
887 vision. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*,
888 pp. 297–306, 2022.
- 889 Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo.
890 Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings*
891 *of the IEEE/CVF International Conference on Computer Vision*, pp. 6023–6032, 2019.
- 892 Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the
893 performance of convolutional neural networks via attention transfer. In *International Conference*
894 *on Learning Representations*, 2017.
- 895 Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. Mixup: Beyond empirical
896 risk minimization. In *International Conference on Learning Representations*, 2018.
- 897 Manyuan Zhang, Guanglu Song, Xiaoyu Shi, Yu Liu, and Hongsheng Li. Three things we need to
898 know about transferring stable diffusion to visual dense prediction tasks. In *European Conference*
899 *on Computer Vision*, pp. 128–145. Springer, 2024.
- 900 Tong Zhang, Pan Ji, Mehrtash Harandi, Richard Hartley, and Ian Reid. Scalable deep k-subspace
901 clustering. In *Computer Vision—ACCV 2018: 14th Asian Conference on Computer Vision, Perth,*
902 *Australia, December 2–6, 2018, Revised Selected Papers, Part V 14*, pp. 466–481. Springer, 2019.
- 903 Zijian Zhang, Zhou Zhao, and Zhijie Lin. Unsupervised representation learning from pre-trained
904 diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 35:22117–
905 22130, 2022.
- 906 Wenliang Zhao, Yongming Rao, Zuyan Liu, Benlin Liu, Jie Zhou, and Jiwen Lu. Unleashing text-
907 to-image diffusion models for visual perception. In *Proceedings of the IEEE/CVF International*
908 *Conference on Computer Vision*, pp. 5729–5739, 2023a.
- 909 Wenliang Zhao, Yongming Rao, Zuyan Liu, Benlin Liu, Jie Zhou, and Jiwen Lu. Unleashing text-
910 to-image diffusion models for visual perception. In *Proceedings of the IEEE/CVF International*
911 *Conference on Computer Vision*, pp. 5729–5739, 2023b.

918 Daquan Zhou, Zhiding Yu, Enze Xie, Chaowei Xiao, Animashree Anandkumar, Jiashi Feng, and
919 Jose M Alvarez. Understanding the robustness in vision transformers. In *International Conference*
920 *on Machine Learning*, pp. 27378–27394. PMLR, 2022.

921
922 Zikai Zhou, Shitong Shao, Lichen Bai, Zhiqiang Xu, Bo Han, and Zeke Xie. Golden noise for
923 diffusion models: A learning framework. *arXiv preprint arXiv:2411.09502*, 2024a.

924 Zikai Zhou, Yunhang Shen, Shitong Shao, Linrui Gong, and Shaohui Lin. Rethinking centered
925 kernel alignment in knowledge distillation. In *Proceedings of the Thirty-Third International Joint*
926 *Conference on Artificial Intelligence*, pp. 5680–5688, 2024b.

927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

972	APPENDIX: TABLE OF CONTENTS	
973		
974		
975	A An overview of the paper	21
976		
977		
978	B Detailed Preliminaries	21
979	B.1 Denoising diffusion probabilistic model (DDPM)	21
980	B.2 Denoising diffusion implicit models (DDIM)	21
981		
982		
983	C Random samples from <i>Tench</i> class generated by DiT	22
984		
985		
986	D Motivation and theoretical explanation of CLARID	22
987		
988	E The layer index for Jacobian calculation	24
989		
990	F The generalization of CLARID	25
991	F.1 Fine-grained control of CanoReps with text conditioning	25
992	F.2 CLARID is compatible with the EDM sampler and UViT architecture	26
993		
994		
995	G Details of the ImageNet20 experiment	27
996	G.1 Selecting the optimal layer and time step for feature extraction	28
997	G.2 Comparison between CLARID combined with CFG and pure CFG	28
998	G.3 On using NMI for measuring feature quality and the effectiveness of <i>CaDistill</i>	29
999		
1000		
1001		
1002	H Choosing hyperparameters t_e and n for CLARID	30
1003	H.1 Finding the optimal t_e	30
1004	H.2 Choosing the total number of extraneous directions n for adaptively selecting k	31
1005	H.2.1 Stable Diffusion 2.1	33
1006	H.3 Failure cases and discussion	33
1007		
1008		
1009		
1010	I Computational costs and scalability of CLARID	36
1011		
1012	J Quantitative analysis of Canonical Samples	37
1013		
1014		
1015	K Details of the proof-of-concept experiment	37
1016	K.1 Data generation process	37
1017	K.2 Architecture and training details	38
1018	K.3 CLARID and baseline configuration	38
1019	K.4 Results and analysis	39
1020		
1021		
1022		
1023	L Training and evaluation details of <i>CaDistill</i>	39
1024	L.1 Adversarial robustness benchmarking	39
1025	L.2 Generalization ability evaluation	40

1026	M More results on ImageNet	41
1027		
1028	M.1 The performance of \mathcal{L}_{dist} alone in feature distillation	41
1029	M.2 <i>CaDistill</i> is effective with different data augmentation strategies	42
1030	M.3 <i>CaDistill</i> improves the student’s black-box adversarial robustness	43
1031		
1032	M.4 Extracting class semantics using class tokens in vision transformer	43
1033	M.5 Details of the Background Challenge	44
1034		
1035	M.6 On the reproduction of RepFusion	44
1036	M.7 Ablation studies	45
1037		
1038	M.7.1 Number of CanoReps	45
1039	M.7.2 The necessity of \mathcal{L}_{align} and \mathcal{L}_{cano} and their balance	45
1040	M.7.3 The necessity of \mathcal{L}_{dist}	45
1041		
1042	M.7.4 The weights of losses, λ_{cs} , λ_{dist} , λ_{cka}	46
1043	M.7.5 The design of \mathcal{L}_{cano}	46
1044		
1045	M.7.6 The magnitude of classifier-free guidance	47
1046	M.8 Generalization of <i>CaDistill</i> to different student architectures	48
1047	M.9 Class-wise robustness improvement	48
1048		
1049	M.10 The variation of extraneous directions within and between classes	49
1050		
1051	N Significance of the quantitative results and discussion	49
1052		
1053	N.1 Significance of the improvements brought by <i>CaDistill</i>	49
1054	N.2 Origin of improvements	49
1055	N.3 Data efficiency of <i>CaDistill</i>	50
1056		
1057	N.4 Difference of the results on CIFAR10 and ImageNet	50
1058		
1059	O Discussion and potential applications of CLARID	51
1060		
1061	P On more sophisticated feature distillation frameworks	52
1062		
1063	Q More visual results	52
1064		
1065	R Broader Impact	52
1066		
1067		
1068	S License information	58
1069		
1070	S.1 Datasets information and license	58
1071	S.2 Model and code license	58
1072		
1073	T Large Language Models usage	58
1074		
1075		
1076		
1077		
1078		
1079		

A AN OVERVIEW OF THE PAPER

We outline the paper in Figure 8.

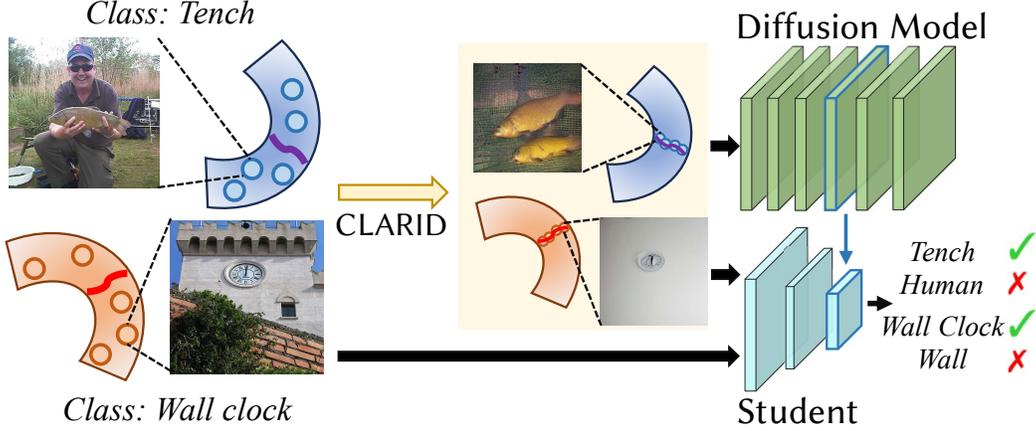


Figure 8: An overview of our paper. We identify the *CanoReps* via **C**anonical **L**atent **R**epresentation **I**dentifier (CLARID) inside conditional diffusion models (CDMs) as a compact family of latent codes that contain the core class information with minimal class-irrelevant signals. These prototypes lie on low-dimensional manifolds (red \sim and purple \sim tildes) within the latent space of CDMs. Leveraging the *CanoReps*, we design a diffusion-based **f**eature **d**istillation paradigm, improving the adversarial robustness and generalization of downstream models in image classification.

B DETAILED PRELIMINARIES

B.1 DENOISING DIFFUSION PROBABILISTIC MODEL (DDPM)

DDPM (Ho et al., 2020) models the generation process as an inversion of a fixed forward Gaussian diffusion $q(\mathbf{x}_{1:T} | \mathbf{x}_0) := \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})$. The forward kernel $q(\mathbf{x}_t | \mathbf{x}_{t-1})$ is described in Eq. 7.

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) =: \mathcal{N}\left(\sqrt{\frac{\alpha_t}{\alpha_{t-1}}} \mathbf{x}_{t-1}, \left(1 - \frac{\alpha_t}{\alpha_{t-1}}\right) \mathbf{I}\right), \quad (7)$$

where $\{\beta_t\}_{t=1}^T$ is the variance schedule, $\alpha_t = \prod_{k=1}^t (1 - \beta_k)$. The inversion process is defined as $p_\theta(\mathbf{x}_{0:T}) := p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)$, where $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. θ denotes the parameter set of a trainable noise predictor \mathbf{f}_θ . A single reverse step is formalized in Eq. 8.

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{1 - \beta_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{\alpha_t}} \mathbf{f}_{\theta,t}(\mathbf{x}_t) \right) + \sqrt{\beta_t} \boldsymbol{\epsilon}_t, \quad \boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (8)$$

$\mathbf{f}_{\theta,t}$ means that the noise predictor receives t as a conditional input.

B.2 DENOISING DIFFUSION IMPLICIT MODELS (DDIM)

DDIM (Song et al., 2021a) proposes a non-Markovian forward diffusion process, implying the parametrization described in Eq. 9.

$$q_\xi(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\sqrt{\alpha_{t-1}} \mathbf{x}_0 + \sqrt{1 - \alpha_{t-1} - \xi_t^2} \frac{\mathbf{x}_t - \sqrt{\alpha_t} \mathbf{x}_0}{\sqrt{1 - \alpha_t}}, \xi_t^2 \mathbf{I}\right), \quad (9)$$

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149



1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162

Figure 9: Random samples generated by DiT 256×256 (Peebles & Xie, 2023) when conditioned on the *Tench* class (0th class in ImageNet).

where $\xi_t = \eta \sqrt{\frac{1-\alpha_{t-1}}{1-\alpha_t}} \sqrt{1 - \frac{\alpha_t}{\alpha_{t-1}}}$. The reverse step is described in Eq. 10.

$$\mathbf{x}_{t-1} = \left(\frac{\mathbf{x}_t - \sqrt{1-\alpha_t} \mathbf{f}_\theta(\mathbf{x}_t)}{\sqrt{\alpha_t}} \right) + \sqrt{1-\alpha_{t-1} - \xi_t^2} \mathbf{f}_\theta(\mathbf{x}_t) + \xi_t \epsilon_t. \quad (10)$$

DDIM and other parametrizations of the diffusion process (Karras et al., 2022; 2024b) can perform inversion on the input sample, retaining certain semantic information of it at $t = T$ (Zhou et al., 2024a).

C RANDOM SAMPLES FROM *Tench* CLASS GENERATED BY DiT

We sample random images using DiT-XL 256×256 (Peebles & Xie, 2023) from the *Tench* class in ImageNet. We use a classifier-free guidance scale of 1.5, which is the one used in the original paper that achieves the best generation quality. The results are shown in Figure 9. Most images that we observe contain an angler.

D MOTIVATION AND THEORETICAL EXPLANATION OF CLARID

Motivation. We show that the extraneous directions (Section 3.2) carry semantics that are not related to the class identity in Figure 10. We adopt the method proposed by Park et al. (2023a). The editing focuses mostly on the background and preserves the class identity, as long as the movement does not orthogonalize the latent code and the editing direction. This phenomenon motivates our experiments.

Theoretical explanation. We provide an explanation of CLARID from the perspective of the change of mutual information. We denote X as the image, Y as the corresponding label, and Z as the CDM feature at time step t_e . Because the CDM is trained to model the full $p(X|Y)$, its internal feature Z will contain a part of the information that is predictive of Y , and the remaining variability in X that is largely independent of Y , such as background, co-occurring objects, and style. We have the following decomposition of the mutual information $I(Z; X, Y)$:

$$I(Z; X, Y) = I(Z; X) + I(Z; Y|X) = I(Z; Y) + I(Z; X|Y)$$

We can get $I(Z; X) = I(Z; Y) + I(Z; X|Y) - I(Z; Y|X)$. $I(Z; Y)$ measures how much of the latent is aligned with class semantics. $I(Z; X|Y)$ represents how much extra information Z contains about the specific X given Y . $I(Z; Y|X)$ is zero, assuming that the label is correct and unique (and hence $I(Z; Y|X) \leq H(Y|X) = 0$) in the dataset. CLARID projects the full Z onto the orthogonal complement of the top Jacobian singular directions, as introduced in Section 3.2. By

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

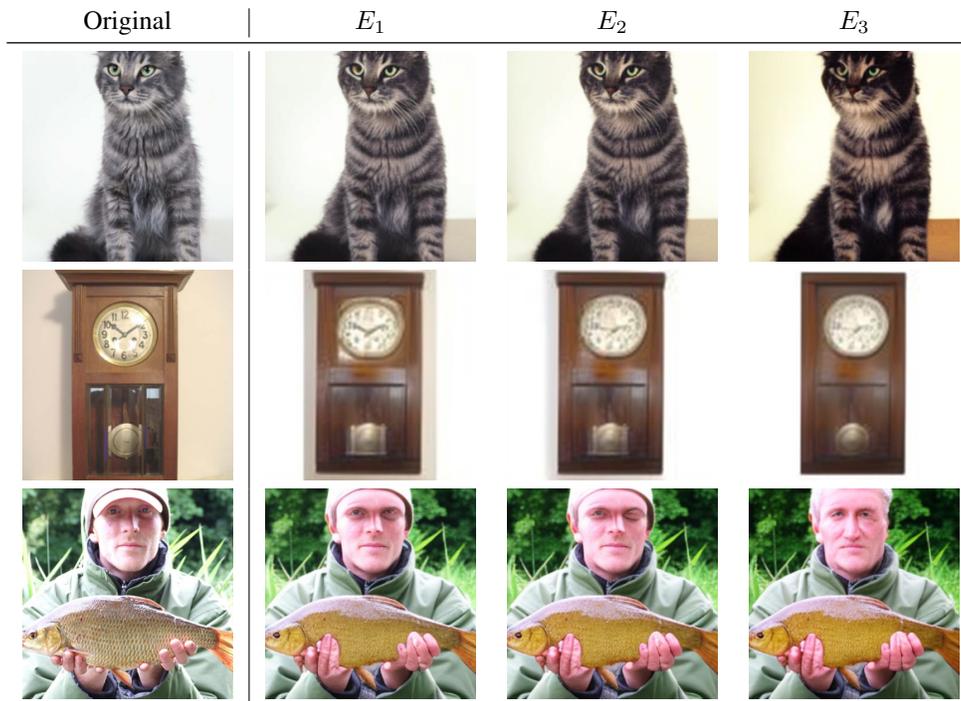


Figure 10: Moving along extraneous directions (Section 3.2) will alter the appearance of the image while preserving the class identity. E_i represents the editing strength.

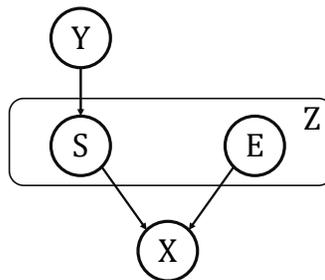


Figure 11: A toy generative model. X is the generated data. Y denotes the label of X . The latent Z is partitioned into two parts S and E , where S depends on Y while E is independent of the label. Given Y , the generative model will have to encode all the variance of X into E to successfully model the full $p(X|Y)$. CLARID identifies those high-variance components given Y and remove them from Z to preserve core class information S , creating an information bottleneck without retraining the CDM.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

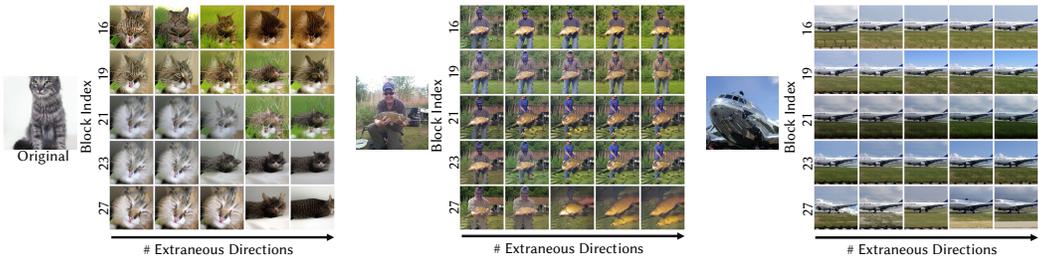


Figure 12: Canonical Samples when computing extraneous directions at different layers in a DiT (Peebles & Xie, 2023), with $t_e = 0.8T$. Note that we do not use CFG after the extraneous directions projection, to present a more straightforward comparison between different layers. We choose $l = 27$, *i.e.* the last layer, to ensure an adequate change in the input.

construction and by our editing experiments in Figure 10, moving along those directions changes the image appearance while keeping the class identity unchanged, so they are a proxy for directions that contribute significantly to $I(Z; Y)$ but not to $I(Z; X|Y)$. Denote the new latent obtained by CLARID as \tilde{Z} . CLARID reduces the total information that Z contains about X , as it removes certain things from Z , hence $I(\tilde{Z}; X) < I(Z; X)$. Meanwhile, it aims to keep the class-relevant term $I(\tilde{Z}; Y) \approx I(Z; Y)$ as unchanged as possible. From this perspective, CLARID implicitly creates an internal information bottleneck inside a pretrained CDM. Therefore, the fraction of latent information that is label-relevant is increased by CLARID. As shown in Figure 3 and 30 in our submitted version, after removing the extraneous directions, the normalized mutual information between the feature cluster assignments and the ground truth labels increases, indicating that the remaining latent is more predictive of the semantic label Y .

We develop a toy generative model as shown in Figure 11. Consider that the full distribution $p(X|Y)$ can be decomposed into $p(Y)p(Z|Y)p(X|Z)$. The latent Z can be further decomposed as $Z = S + E$, where S depends on the label Y and E is independent of it. S is known when Y is given. Therefore, to model the full $p(X|Y)$, the generative model must encode all variance of X in E . CLARID identifies those high-variance components given Y and removes them from Z to preserve core class information S . After removing E , $I(Z; X)$ decreases to $I(S; X)$ while $I(Z; Y)$ does not change and equals $I(S; Y)$, hence the information bottleneck effect.

This perspective is related to the information-theoretic goals of β -TCVAE (Chen et al., 2018) and DBAE (Kim et al.), but differs in where and how the information split happens. In β -TCVAE, the focus is on the total correlation (TC) term in the ELBO decomposition. Increasing the weight on TC explicitly penalizes entangled latent dimensions while preserving information that is predictive of the ground truth factors. Conceptually, the extraneous directions that CLARID identifies play a similar role. They correspond to high-variance directions that lead to many visually different generations within the same class, and thus tend to contribute more to $I(Z; X|Y)$ than to $I(Z; Y)$. In this sense, the effect of CLARID (from our information bottleneck explanation) parallels that of the TC penalty. Crucially, however, CLARID implements this bottleneck in a pretrained CDM, whereas β -TCVAE and DBAE achieve the information constraints by modifying the training objective, adding new modules, and retraining the entire model.

E THE LAYER INDEX FOR JACOBIAN CALCULATION

In Section 3.2, we treat the CDM as a feature extractor for calculating the Jacobian. We visualize the effect of selecting different l , *i.e.* the layer index for computing the Jacobian, in Figure 12 with t_e fixed to $0.8T$, and Figure 13 with k chosen by the method described in Section 3.2.1. While in some cases different layers can yield similar effects, only the top one, *i.e.* $l = 27$, can ensure an adequate change in the output image, or the background can remain unchanged in certain cases.

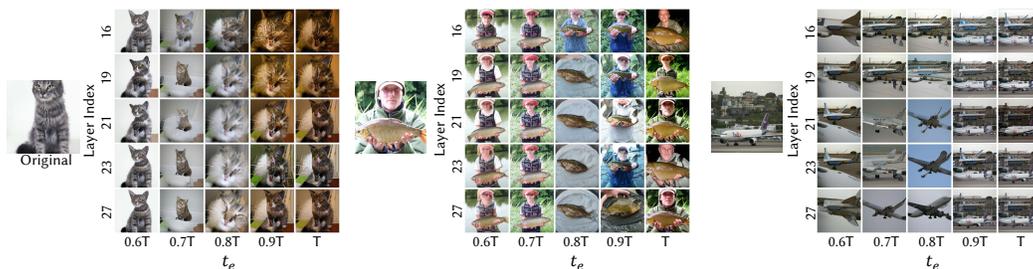


Figure 13: Canonical Samples when computing extraneous directions at different layers in a DiT (Peebles & Xie, 2023). For each image, we choose the number of extraneous directions to be removed automatically according to the method in Section 3.2.1. Note that we do not use CFG after the extraneous directions projection, to present a more straightforward comparison between different layers. We choose $l = 27$, i.e. the last layer, to ensure an adequate change in the input.

For Stable Diffusion (Rombach et al., 2022), we follow the practice in Park et al. (2023a) to select the layer index. Concretely, we extract the features after the middle block to ensure that the extraneous directions are semantically meaningful (Jeong et al., 2024; Kwon et al., 2023). We adopt the same strategy for all UNet-based CDMs, including the one used in our CIFAR10 experiments in Section 4.1 and the CDM in the EDM framework (Karras et al., 2022; 2024b).

F THE GENERALIZATION OF CLARID

F.1 FINE-GRAINED CONTROL OF CANOREPS WITH TEXT CONDITIONING

CLARID naturally extends to text-conditioning CDMs. Text-conditioning offers a more flexible control over where CanoReps lie than one-hot label conditioning. We show visual results in Figure 14 and 15. The used CDM, a Stable Diffusion 2.1 (Rombach et al., 2022), successfully adapts to different text prompts on the same image, which is in line with previous findings (Park et al., 2023a). In Figure 15, the CDM finds a CanoRep that does not exist in the real world, but all the components in it are real, e.g. the water and the airplane. The results demonstrate that it is possible to perform a more fine-grained control over where CanoReps lie. We believe investigating the relationship between the CanoReps and the text conditioning on the same image can reveal the image understanding capability of CDMs, which we leave as a promising future direction.

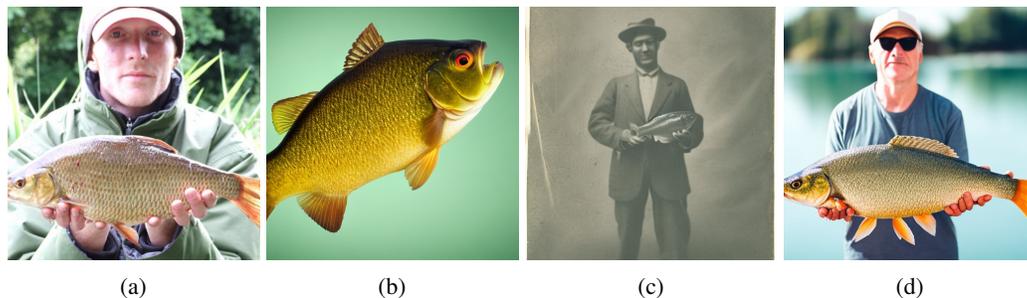
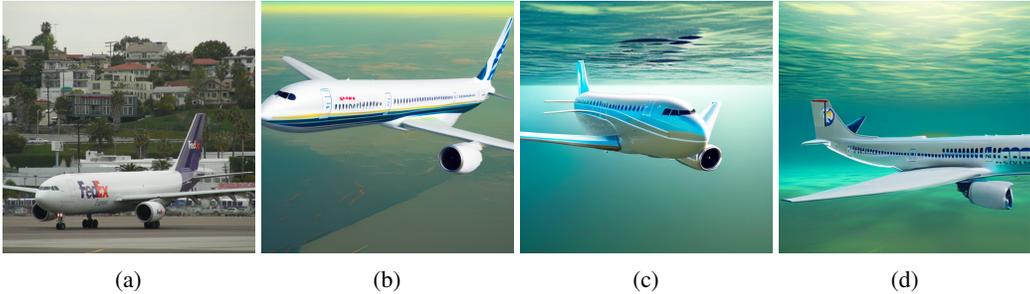


Figure 14: (a): The original image from the class *Tench*. (b): The Canonical Sample obtained with prompt: *a photo of tench*. (c): The Canonical Sample obtained with prompt: *a photo of a man holding a fish*. (d): An image generated with CFG using the same prompt as in (c), using the same starting noise as in (b) and (c).

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359



1360 Figure 15: (a): The original image from the class *Airliner*. (b): The Canonical Sample obtained with
1361 prompt: *a photo of airliner*. (c): The Canonical Sample obtained with prompt: *an airplane flying*
1362 *under water*. (d): An image generated with CFG using the same prompt as in (c), using the same
1363 starting noise as in (b) and (c).

1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379

	Class	Original	CLARID	Class	Original	CLARID
EDM (Karras et al., 2022; 2024b)	Tench			Wine Bottle		
					Screen	

1380 Figure 16: Preliminary results on Canonical Samples generated by the EDM framework (Karras
1381 et al., 2022; 2024b). The left column is an EDM trained on ImageNet 64×64 , the right one is on
1382 ImageNet 512×512 . We implement the inversion sampler of EDM. It indicates that the main idea
1383 behind CLARID is also effective when facing inputs with different resolutions. Note that we do not
1384 use CFG or Autoguidance (Karras et al., 2024a) to improve the visual quality, to provide a more
1385 straightforward insight into the effectiveness of CLARID on EDM.

1386
1387
1388
1389

F.2 CLARID IS COMPATIBLE WITH THE EDM SAMPLER AND UViT ARCHITECTURE

1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

The main idea behind CLARID is to identify the latent vectors that carry non-discriminative information and render the latent code of a F_{inv} -inverted sample orthogonal to them. Such a formulation does not require any knowledge about the sampler or the architecture of the model, as long as the DM has enough capability to model the conditional data distribution. Here, we demonstrate preliminary results on the generalization of CLARID. Specifically, we test the main idea behind CLARID on the EDM sampler (Karras et al., 2022; 2024b) with a UNet-based CDM, and a UViT (Bao et al., 2023) model using the same DDIM sampler as in the main paper. All CDMs are trained on ImageNet. We implement the inversion sampler of EDM. Here, we aim at showing the effectiveness of the identification of extraneous directions but not on the feature quality. Hence, we do not perform the same analysis as in Section 3.2.1. We present some visual results to demonstrate the intuitive summary of the categorical semantics, as shown in Figure 16 and 17. We believe investigating the relationship between the performance of the DM in generative tasks and it as teacher in *CaDistill* is promising, as previous works have found (Wang et al., 2023; Xiang et al., 2023), and leave it as a future work.

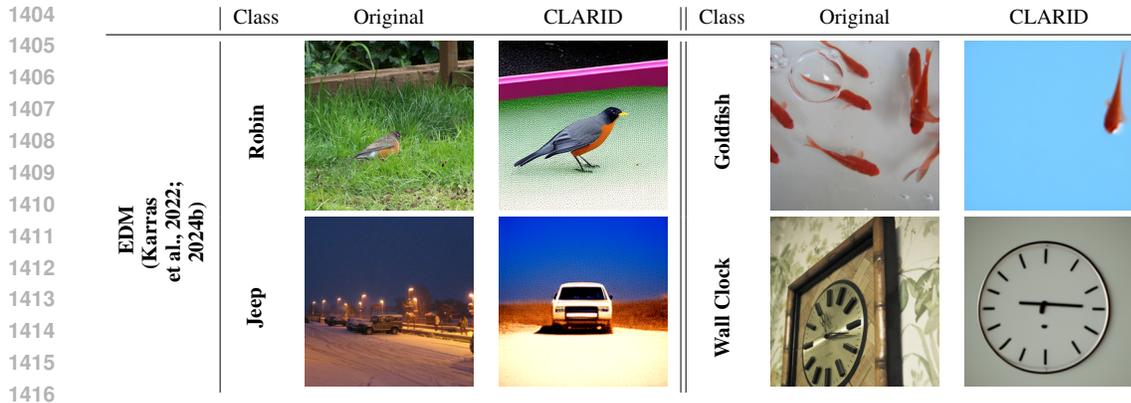


Figure 17: Preliminary results on Canonical Samples generated by a ImageNet-256 UViT (Bao et al., 2023) using DDIM (Song et al., 2021a) sampler. It indicates that the main idea behind CLARID is effective on different architectures of CDMs. Note that we do not use CFG to improve the visual quality, to provide a more straightforward insight into the effectiveness of CLARID on UViT.

G DETAILS OF THE IMAGENET20 EXPERIMENT

To develop the CLARID framework, we perform quantitative experiments on a 20-class subset of ImageNet (Deng et al., 2009). This choice balances between computational costs and statistical significance. Our goal is to construct a subset with a clear structure while keeping the class-separation task neither too easy nor too hard. On the one hand, the classes should be dissimilar enough so that separation is meaningful; on the other hand, if they are too dissimilar, separation becomes trivial, and it is hard to draw reliable conclusions. To balance this, we start from the widely used 16-way ImageNet split (Geirhos et al., 2018b;a; Subramanian et al., 2024; Gavrikov & Keuper, 2024), and randomly select one class from each of the 16 superclasses. We then increase the difficulty by randomly selecting 4 additional classes from the bird and dog superclasses, which contain the most subclasses. The resulting subset has a similar number of classes from the two main partitions of ImageNet, animals and artifacts, ensuring a well spread over ImageNet1K.

To demonstrate the robustness and generalization of our method, we conduct three runs by changing the selection of the 20 classes and report standard deviation in our main paper Figure 3. To ensure reproducibility, we list all the classes in the three runs:

- Run1. Indices of the 20 classes: [15, 95, 146, 151, 211, 242, 281, 294, 385, 404, 407, 409, 440, 444, 499, 544, 579, 717, 765, 814]. The corresponding class names: ['robin', 'jacamar', 'albatross', 'Chihuahua', 'vizsla', 'boxer', 'tabby', 'brown bear', 'Indian elephant', 'airliner', 'ambulance', 'analog clock', 'beer bottle', 'bicycle-built-for-two', 'cleaver', 'Dutch oven', 'grand piano', 'pickup', 'rocking chair', 'speedboat'].
- Run2. [7, 94, 97, 143, 152, 266, 281, 294, 385, 405, 409, 436, 440, 499, 554, 555, 559, 671, 687, 766]. The corresponding class names: ['cock', 'fireboat', 'drake', 'oystercatcher', 'toy poodle', 'Irish setter', 'folding chair', 'brown bear', 'Indian elephant', 'airship', 'Japanese spaniel', 'tabby', 'beach wagon', 'cleaver', 'beer bottle', 'fire engine', 'analog clock', 'mountain bike', 'organ', 'rotisserie'].
- Run3. [15, 96, 146, 152, 212, 268, 282, 295, 386, 405, 436, 530, 623, 671, 687, 717, 720, 765, 766, 814]. The corresponding class names: ['robin', 'toucan', 'albatross', 'Japanese spaniel', 'English setter', 'Mexican hairless', 'tiger cat', 'American black bear', 'African elephant', 'airship', 'beach wagon', 'digital clock', 'letter opener', 'mountain bike', 'organ', 'pickup', 'pill bottle', 'rocking chair', 'rotisserie', 'speedboat'].

We plot the pair-wise Wu-Palmer (WUP) distances (Wu & Palmer, 1994) between the selected classes in Figure 18. The selected classes can be similar or dissimilar to each other, demonstrating certain structures. This is appropriate for analyzing the class separability in our case. We choose 50 images from the ImageNet20, building a 1000-sample dataset for the following analysis.

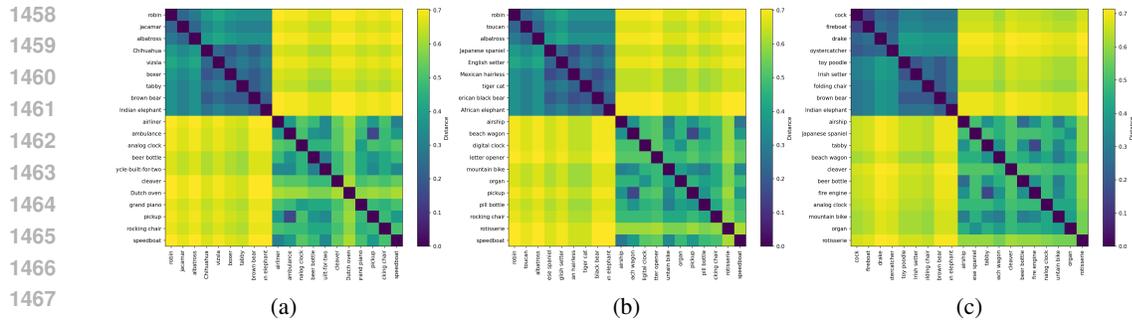


Figure 18: Wu-Palmer (WUP) distances (1.0-WUP similarity (Wu & Palmer, 1994)) between the classes in ImageNet20 across three runs. (a): Run1. (b): Run2. (c): Run3. The class relationships are structured, hence appropriate for analyzing the class separability in our case. **Note that we select the classes from the 16-way ImageNet (Geirhos et al., 2018b;a; Subramanian et al., 2024; Gavrikov & Keuper, 2024) to ensure a fair task difficulty, which can lead to a few class overlaps between runs.**

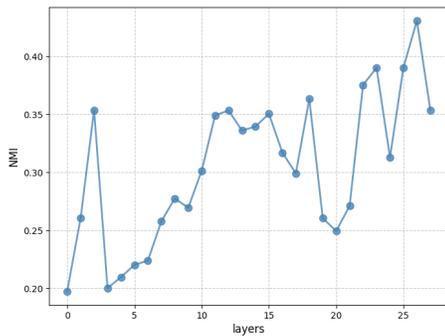


Figure 19: NMI v.s. layers on ImageNet20 in a DiT (Peebles & Xie, 2023), fixing $t_r = 0.1T$. We choose the penultimate layer, *i.e.* the 27th layer (the figure has a 0-start index), in all our experiments.

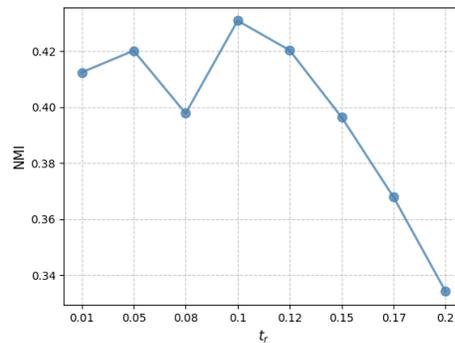


Figure 20: NMI v.s. feature extraction time step (t_r) on ImageNet20 using a DiT (Peebles & Xie, 2023), fixing the layer index to be 27. We choose $t_r = 0.1T$ in all our experiments.

G.1 SELECTING THE OPTIMAL LAYER AND TIME STEP FOR FEATURE EXTRACTION

We choose an ImageNet 256×256 -pretrained DiT-XL (Peebles & Xie, 2023) model for the quantitative analysis. We consider the outputs of all 28 ViT blocks in it. For the time step, we select $t_r = \{0.01, 0.05, 0.08, 0.1, 0.12, 0.15, 0.17, 0.2\}$. We perform K-means clustering on the feature map after average pooling, and compute the normalized mutual information between the cluster assignments and the ground truth class labels. The average pooling reduces noise in the feature map, making the cluster more accurate and compact. The results for different layers and time steps are shown in Figure 19, 20. Our conclusion on the time step for feature extraction is consistent with previous studies (Mukhopadhyay et al., 2023; Yang & Wang, 2023) that adopt linear probing on the features for quantifying the feature quality, albeit with different DM architectures. Such a result provides **evidence on the validity of our metric, *i.e.* normalized mutual information.**

G.2 COMPARISON BETWEEN CLARID COMBINED WITH CFG AND PURE CFG

In Section 3.4 in the main paper, we show that Canonical Samples, obtained via CLARID combined with CFG, contain fewer class-irrelevant components than samples obtained via CFG. We now show that Canonical Samples are more separable in the CDM feature space than using pure CFG to quantitatively demonstrate that they encode different information. We compare the cluster quality between the features corresponding to the two kinds of samples, namely Canonical Features and

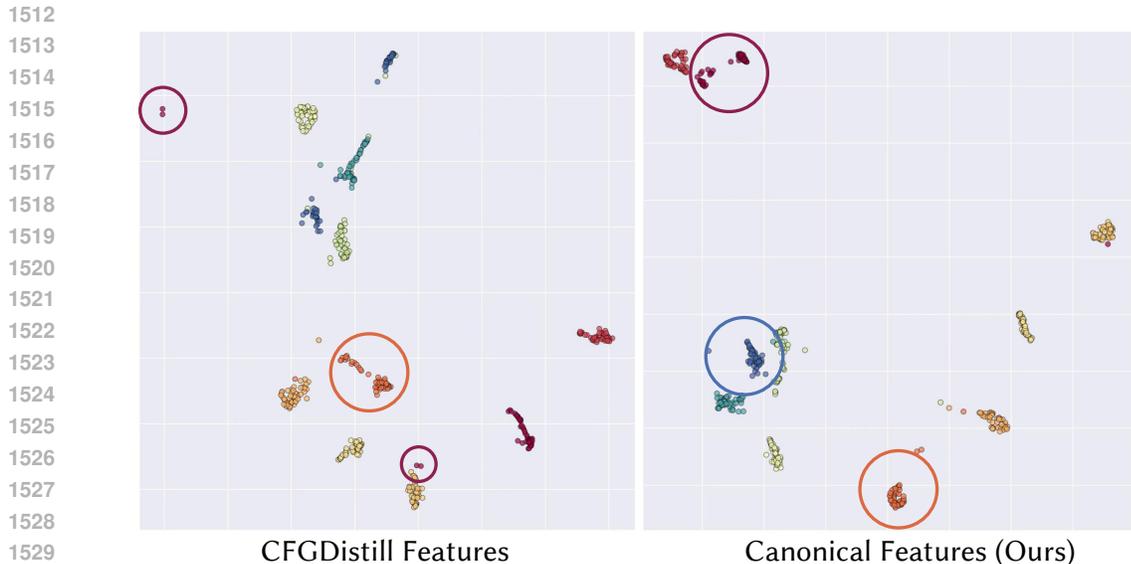


Figure 21: A 2D UMAP (McInnes et al., 2018) projection of the CDM feature space, showing clusters for ten classes. In this case, Canonical Features are produced by CLARID combined with CFG. Colors indicate classes. Canonical Features form tighter, more uniformly shaped clusters, whereas CFGDistill features often result in irregular cluster shapes (orange), more outliers (red), and disconnected regions even for samples sharing the same class label (blue).

CFGDistill features, using the pipeline in Section 3.2.1. Canonical Features achieve an NMI score of 0.7762 while CFGDistill features achieve 0.7308. The UMAP projections of the two types of features are shown in Figure 21. Canonical Features form better separated clusters, implying that they capture different information from using pure CFG.

G.3 ON USING NMI FOR MEASURING FEATURE QUALITY AND THE EFFECTIVENESS OF *CaDistill*

While NMI is valid for self-evaluation of the feature quality within the CLARID framework, it is not valid for a direct comparison between the quality of features obtained via different methods, such as CFG. For example, when using CFG on the same 1000 samples after F_{inv} and extracting features at the same t_r , as in CLARID, it can yield a higher NMI (0.6108 in CLARID v.s. 0.7808 with CFG magnitude being 4.0). Adopting CFG after performing extraneous direction projection can also improve the metric (0.7762 in CLARID with CFG being 3.0). However, NMI only examines the compactness and separability of each cluster, ignoring the low-dimensional manifold structure of the data. For example, a line-shaped manifold and a circle-shaped manifold can yield the same NMI, while they capture different characteristics of the data. In an extreme case where the features are constant for all samples belonging to the same class, the NMI will be 1.0, while the features are not meaningful in this case. **In other words, we do not want the student to just learn a "converging" prior on the features belonging to the same classes, but to learn actual class semantics.** This claim is validated in our extended ablation studies in Section M.7, where the student does not perform well when the CFG magnitude is high. These results undermine the notion that the student only mimics the CDM’s label-conditioning embeddings: those embeddings behave as a trivial collapsing prior rather than conveying the encoded semantic structure. A promising future direction is to develop new feature quality metrics or adapt existing ones into CLARID that consider the structure of the data. This metric is also invalid on datasets with simple data structures, such as CIFAR10 (Krizhevsky et al., 2009). The CDM features of the original samples are already perfectly separable and achieve an NMI of 1.0. However, as shown in Table 1 in the main paper, CLARID can still find more semantically meaningful samples than the original ones, and convey the knowledge to the student via *CaDistill*. While the simple structure of the data prevents the utilization of NMI, it is significantly easier to process than real image datasets. Either calculating the extraneous directions or

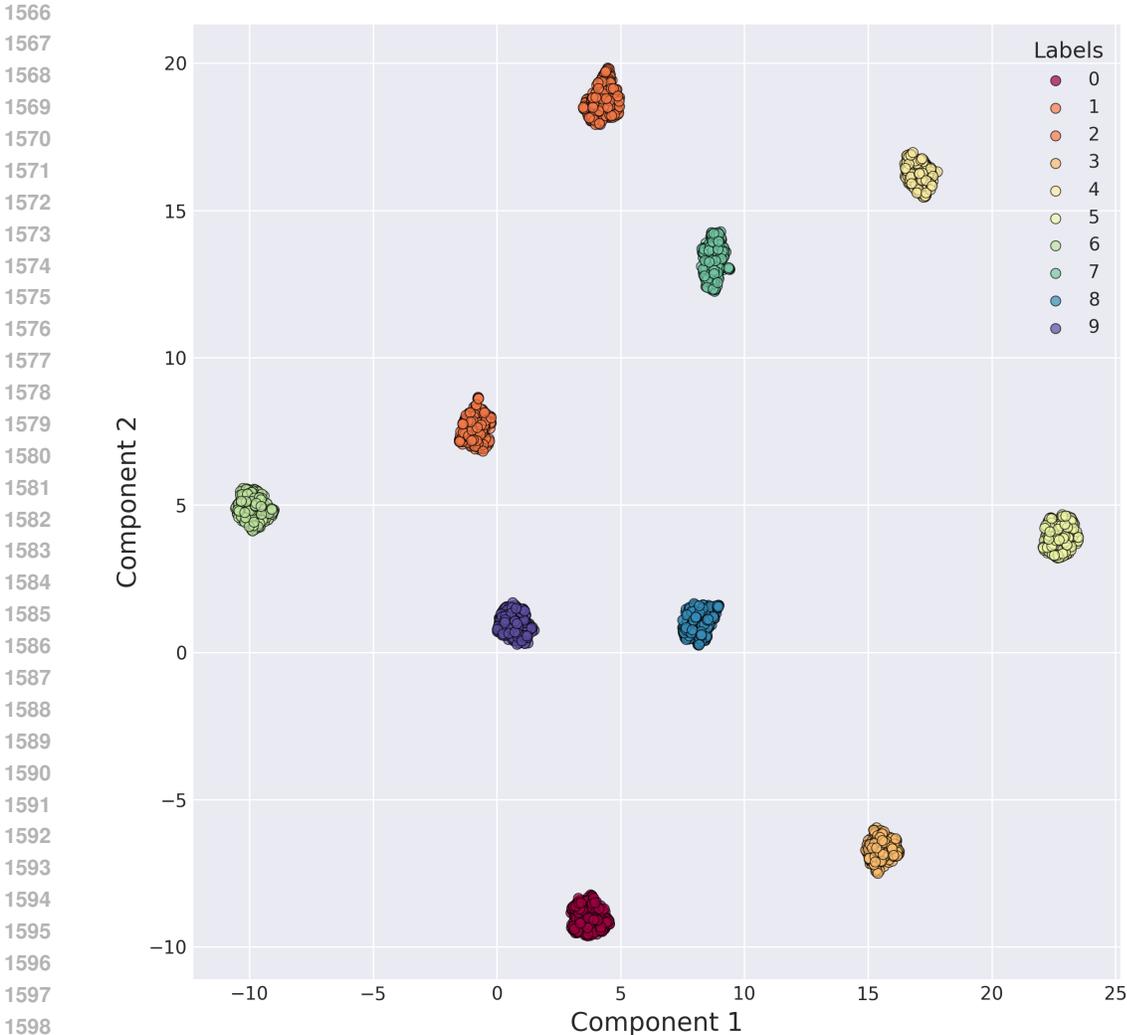


Figure 22: The CDM features of the original samples in CIFAR10 (Krizhevsky et al., 2009) are already separable and achieve an NMI of 1.0. However, CLARID is still effective in this case, as shown in Table 1.

training with *CaDistill* requires far less computational resources than large, real-image datasets such as ImageNet (Deng et al., 2009). Therefore, we simply brute-force search the best hyperparameters in CLARID on CIFAR10. The hyperparameters on CIFAR10 are: $t_e = 0.8, n = 10, t_r = 0.13T$, layer = *up.0*. We also only use a 10% subset for obtaining CanoReps, as shown in Table 1.

H CHOOSING HYPERPARAMETERS t_e AND n FOR CLARID

H.1 FINDING THE OPTIMAL t_e

In Section 3.2, we decide t_e by finding the saturation point of classification accuracy on samples generated by our two-stage strategy. We use 3 classifiers and compute the average accuracy of them:

1. A ViT-Large pre-trained on ImageNet12K (Deng et al., 2009). The input size is 224.
2. An ImageNet22k-pre-trained Swin V2 (Liu et al., 2022). The input size is 256.
3. An ImageNet22k-pre-trained ConvNeXt V2 (Woo et al., 2023). The input size is 384.

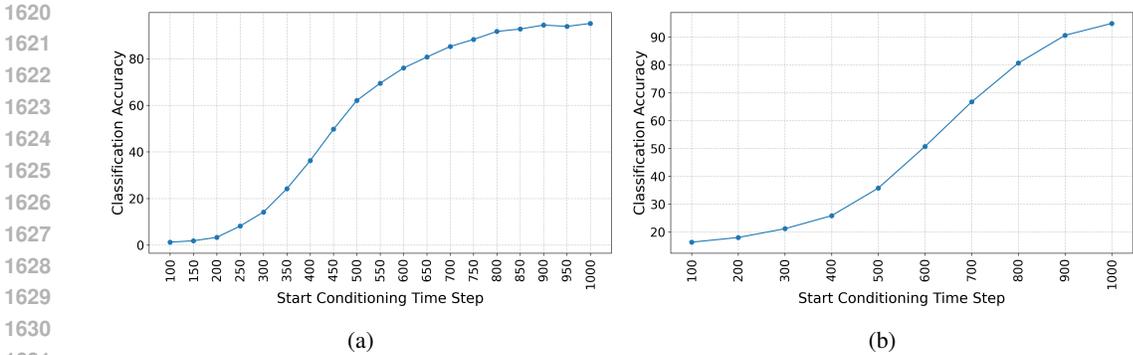


Figure 23: The classification accuracy on samples generated by the two-stage strategy in Section 3.2, using (a) DiT (Peebles & Xie, 2023) and (b) Stable Diffusion (Rombach et al., 2022). The accuracy is averaged over 3 classifiers.

All model weights are downloaded from the PyTorch Image Model library (timm) (Wightman, 2019). The accuracy curve for DiT is shown in Figure 23. The maximum time step T is 1000 for both DiT (Peebles & Xie, 2023) and Stable Diffusion (Rombach et al., 2022). We use DDIM as the sampler and set the total diffusion time step to 50. For each class, we generate 5 images. Hence $m = 5$ in Section 3.2. The prompt template for Stable Diffusion is: *a photo of <class name>*, where the class name is the WordNet name of each ImageNet class (Deng et al., 2009). We identify $t_e = 800$ ($0.8T$) for DiT and $t_e = 1000$ (T) for Stable Diffusion. We provide visual comparisons between the images generated by selecting different t_e in Figure 24. Qualitatively, a small t_e might lead to insufficient changes in the input image, while a large t_e in DiT can fail to make the model aware of the class conditioning, resulting in less meaningful extraneous directions.

H.2 CHOOSING THE TOTAL NUMBER OF EXTRANEIOUS DIRECTIONS n FOR ADAPTIVELY SELECTING k

Our key observation on selecting k , *i.e.* the number of extraneous directions to be projected, is that the effect caused by projecting extraneous directions is not smooth with varying k . That is, projecting one more extraneous directions can cause significant changes in the input. It is because extraneous directions are orthogonal to each other by design, hence there is no guarantee that their semantics are correlated. Importantly, projecting more extraneous directions does not necessarily mean a more separable set of Canonical Features. It can lead to a loss of the class-defining cues, as shown in the fish image in Figure 25b. We show the visual effects of selecting different k in Figure 24, 25. Our method CLARID selects k by adaptively choosing the elbow point on the explained variance ratio (EVR) sequence with total number of extraneous directions being n . The algorithm for finding the elbow point is presented in Algorithm 1.

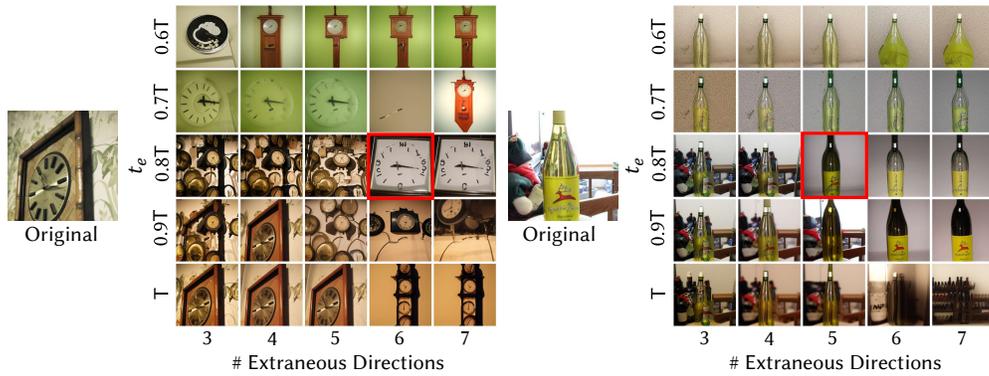
Algorithm 1 Find Elbow via Knee Method(Sequence $S[0 \dots n - 1]$)

```

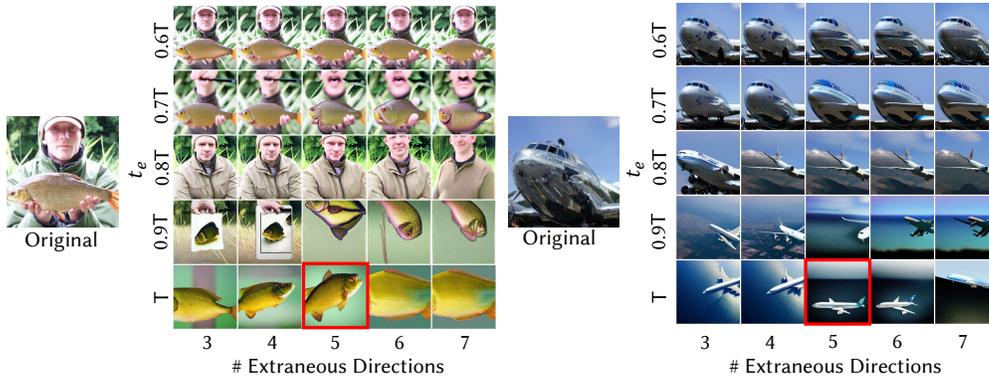
1661 1:  $n \leftarrow |S|$ 
1662 2:  $P_{\text{start}} \leftarrow (0, S[0])$ 
1663 3:  $P_{\text{end}} \leftarrow (n - 1, S[n - 1])$ 
1664 4:  $v \leftarrow P_{\text{end}} - P_{\text{start}}$ 
1665 5:  $u \leftarrow v / \|v\|$ 
1666 6: for  $i = 0$  to  $n - 1$  do
1667 7:    $w \leftarrow (i, S[i]) - P_{\text{start}}$ 
1668 8:    $\text{projLen} \leftarrow w \cdot u$ 
1669 9:    $\text{projVec} \leftarrow \text{projLen} \times u$ 
1670 10:   $\text{perpVec} \leftarrow w - \text{projVec}$ 
1671 11:   $d[i] \leftarrow \|\text{perpVec}\|$ 
1672 12: end for
1673 13:  $k \leftarrow \arg \max_{0 \leq i < n} d[i]$ 
1674 14: return  $k$ 

```

1674
 1675
 1676
 1677
 1678
 1679
 1680
 1681
 1682
 1683
 1684
 1685
 1686
 1687
 1688
 1689
 1690
 1691
 1692
 1693
 1694
 1695
 1696
 1697
 1698
 1699
 1700
 1701
 1702
 1703
 1704
 1705
 1706
 1707
 1708
 1709
 1710
 1711
 1712
 1713
 1714
 1715
 1716
 1717
 1718
 1719
 1720
 1721
 1722
 1723
 1724
 1725
 1726
 1727

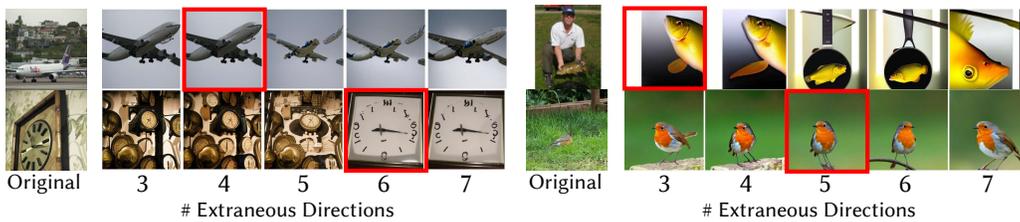


(a) Results of DiT.



(b) Results of Stable Diffusion.

Figure 24: Visual comparisons between the images generated by selecting different t_e and k in CLARID. Red boxes indicate the one automatically selected by our method.



(a) Results of DiT without CFG for a more straightforward comparison. (b) Results of Stable Diffusion with CFG magnitude being 7.5.

Figure 25: Visual comparisons between the images generated by selecting different k in CLARID. Red boxes indicate the one automatically selected by our method.

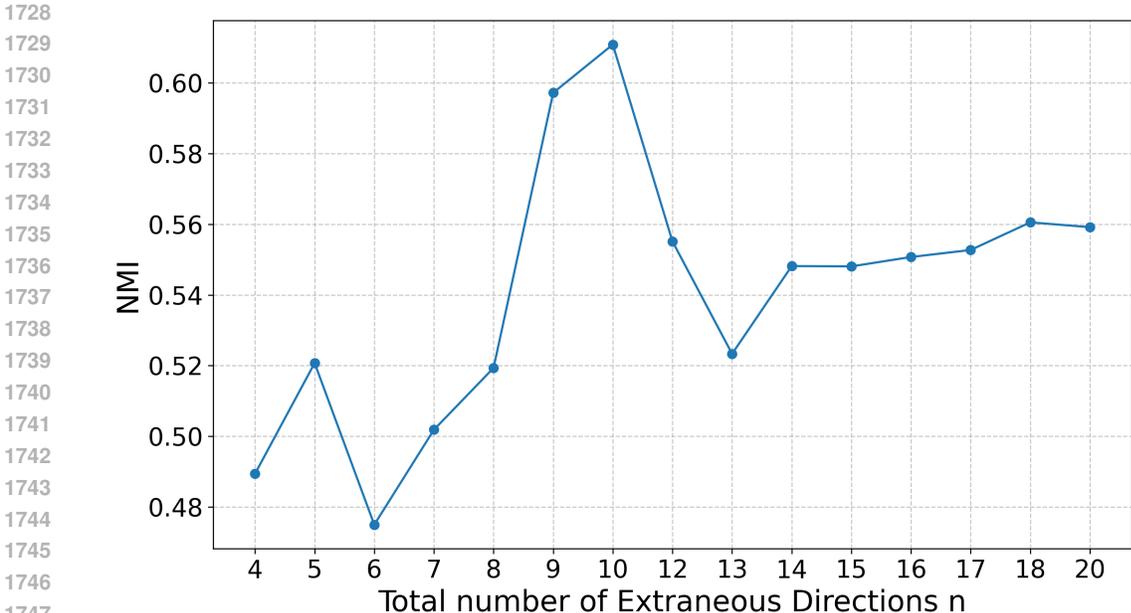


Figure 26: NMI on ImageNet20 v.s. total number of extraneous directions n on DiT. A large n can diminish the discriminative power of the input images, whereas a small one cannot change the inputs too much. Neither case is desired. Hence, we choose $n = 10$ for DiT in our experiment. Note that this is a self-evaluation within the CLARID framework, hence NMI is valid in this case.

We show a quantitative comparison between different n in Figure 26 on ImageNet20. Too large n tends to select a larger k , leading to eliminating too many components in the input image and diminishing discriminative power. A small n will lead to a small k and cannot change the inputs too much. We select $n = 10$ for DiT, and $n = 10$ for SD in our experiment. Note that this n is tailored to specific CDMs, and the same n for DiT and SD is merely a coincidence. We plot the histogram of k when fixing $n = 10$ on 78000 images from ImageNet100, as we used in our experiments in Section M.7, in Figure 27. The selection process does not converge to a single k , supporting the effectiveness of our method. CLARID has certain fault tolerance capacity, *i.e.* slightly changing the number of projected extraneous directions can still result in desired images. For example, in Figure 25, selecting $k = 3$ or $k = 4$ for the airplane image on DiT, or selecting $k = 3 \sim 7$ for the bird image on Stable Diffusion, can all result in desired outputs.

H.2.1 STABLE DIFFUSION 2.1

We provide additional results on Stable Diffusion 2.1 model (Rombach et al., 2022). We first choose $t_r = 0.13T$ and the layer to be up_blocks.1, according to the results in Figure 28, 29. We then perform the same analysis as done in Section 3.2.1 in the main paper and Section H. Figure 30 shows the results. Observe that $t_e = T$, which is $t_e = 0.999$ in the figure, yields the best results when performing fixed- k projection. This validates our saturation-point-based t_e selection, as described in Section 3.2.1 in the main paper. Again, our CLARID produces the highest NMI among all other methods. In Figure 31, we show the different NMI results achieved by selecting different n for deciding k adaptively, and choose $n = 10$.

H.3 FAILURE CASES AND DISCUSSION

We identify some promising future directions and discuss the failure cases. First, we fix t_e for all samples, which can be suboptimal on certain images. We only perform CLARID in a single time step instead of selecting multiple t_e . A series of projecting away extraneous directions has the potential of discarding more class-irrelevant information. Regarding the number of extraneous directions, we only perform experiments on cumulative projection, *i.e.* projecting away the top- k extraneous

1782
 1783
 1784
 1785
 1786
 1787
 1788
 1789
 1790
 1791
 1792
 1793
 1794
 1795
 1796
 1797
 1798
 1799
 1800
 1801
 1802
 1803
 1804
 1805
 1806
 1807
 1808
 1809
 1810
 1811
 1812
 1813
 1814
 1815
 1816
 1817
 1818
 1819
 1820
 1821
 1822
 1823
 1824
 1825
 1826
 1827
 1828
 1829
 1830
 1831
 1832
 1833
 1834
 1835

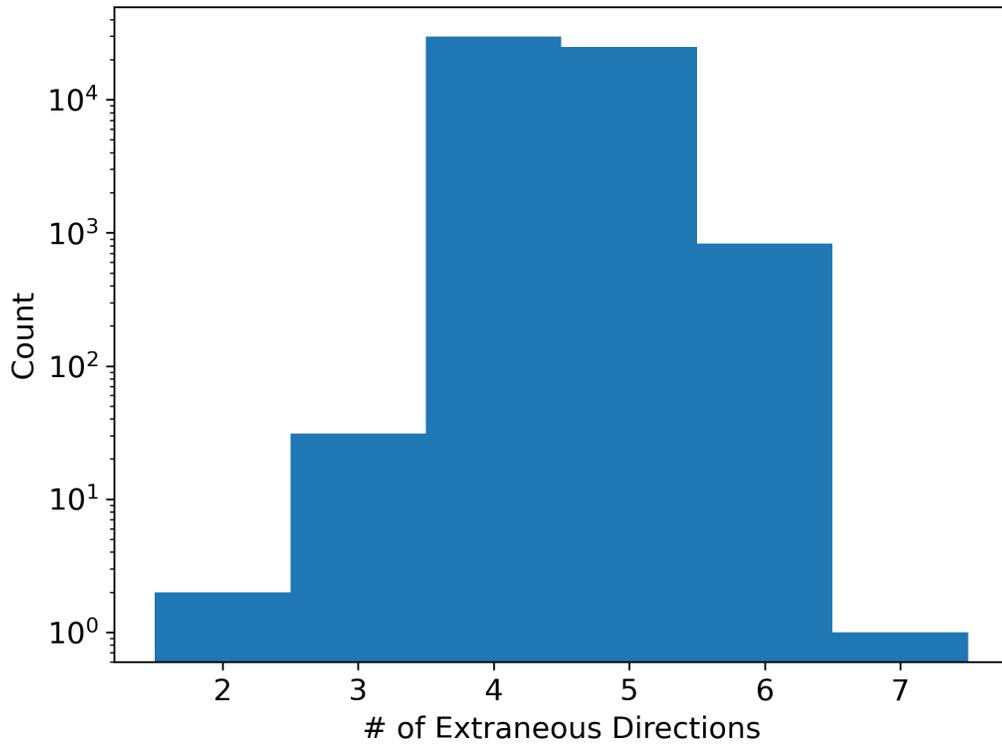


Figure 27: Histogram of k when fixing $n = 10$ on 78000 images from ImageNet100. We use this dataset in our experiments in Section M.7.

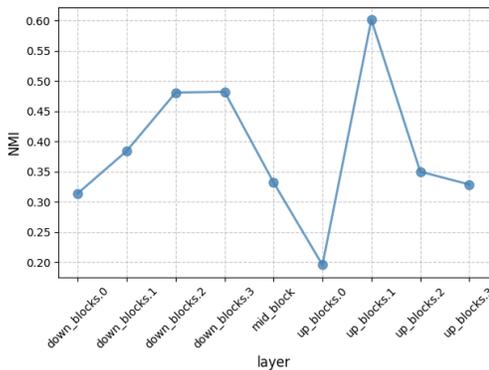


Figure 28: NMI v.s. layers on ImageNet20 in a Stable Diffusion 2.1 (Rombach et al., 2022), fixing $t_r = 0.13T$. We choose the up_blocks.1 layer in all our experiments.

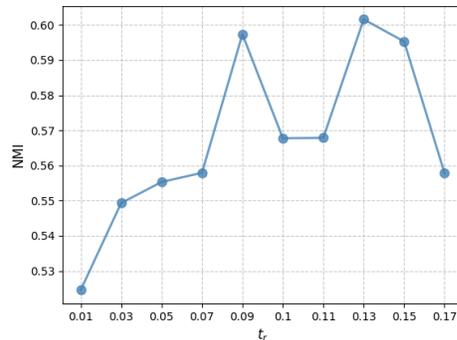


Figure 29: NMI v.s. feature extraction time step (t_r) on ImageNet20 using a Stable Diffusion 2.1 (Rombach et al., 2022), fixing the layer index to be up_blocks.1.

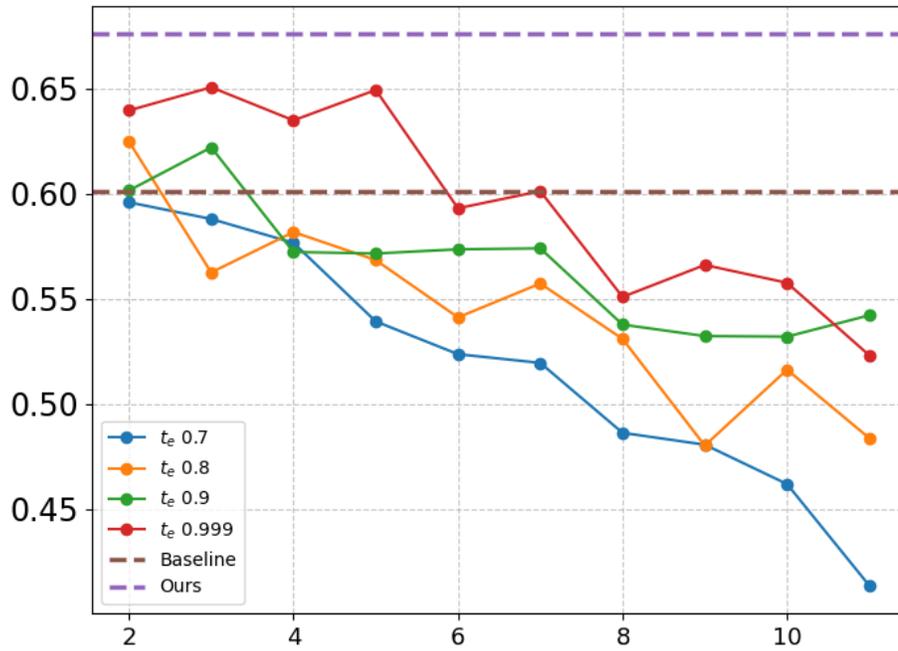


Figure 30: The normalized mutual information (NMI, higher is better) between cluster assignments of Stable Diffusion (SD) features using a Stable Diffusion 2.1 (Rombach et al., 2022) and the ground truth labels. CLARID achieves the highest NMI. Baseline is the original SD features.

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910

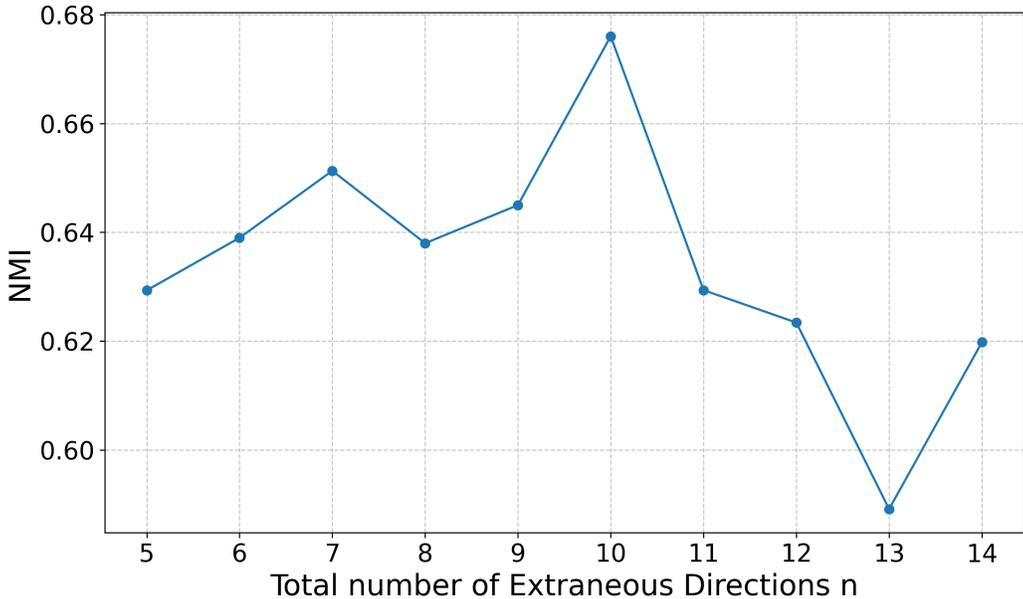


Figure 31: NMI on ImageNet20 v.s. total number of extraneous directions n on a Stable Diffusion 2.1 (Rombach et al., 2022). A large n can diminish the discriminative power of the input images, whereas a small one cannot change the inputs too much. Neither case is desired. Hence, we choose $n = 10$ for Stable Diffusion in our experiment. Note that this is a self-evaluation within the CLARID framework, hence NMI is valid in this case.

1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921

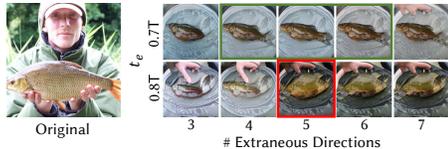


Figure 32: A Failure case of CLARID on selecting t_e . Green boxes are the optimal choice, qualitatively. Red boxes are the ones CLARID selects.

1922
1923
1924
1925
1926
1927
1928
1929
1930
1931

directions. A careful selection of the extraneous directions can contribute to better results. Due to the limits in computational resources, we leave them as future work. Occasionally, CLARID can select suboptimal k , leading to artefacts in the generated images. We show failure cases for t_e and k in Figure 32, 33. The aforementioned future directions can potentially serve as solutions to those failure modes.

1932
1933
1934

I COMPUTATIONAL COSTS AND SCALABILITY OF CLARID

1935
1936
1937
1938
1939
1940
1941
1942
1943

We approximate only the top singular vectors of the Jacobian via the established Jacobian subspace iteration method (Haas et al., 2024; Park et al., 2023a), avoiding a full SVD. Complexity scales quadratically in the number of vectors. We use $n = 10$ per image. Measured on a single Nvidia A100, the runtime is 9.5s for DiT and 22s for Stable Diffusion. A simple speed-up is early truncation: reducing max iters from 100 to 20 and relaxing the threshold from $1e-4$ to $1e-3$ yields 1.5 s (DiT) and 4.4 s (SD) with visually similar outputs. This is valid because we actually use only the top singular vectors, while the remaining ones serve merely to locate the elbow. Developing faster Jacobian singular-vector estimators is a promising future work. On scalability, CLARID is an offline preprocessing step. A single trial suffices for a dataset. In classification (Section M.7), processing only 10% of data attains competitive downstream performance. CLARID works across different

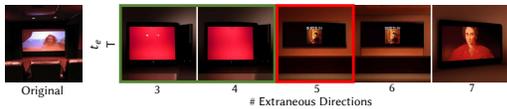


Figure 33: A Failure case of CLARID on selecting k when $n = 10$. Green boxes are the optimal choice, qualitatively. Red boxes are the ones CLARID selects.

Table 3: Quantitative comparison between Training samples, the samples used in CFGDistill (CFGDistill samples), and Canonical Samples. All metrics are computed following the pipeline in (Dhariwal & Nichol, 2021).

	FID	sFID	Inception Score	Precision	Recall	Classification Accuracy
Training samples	0.8	3.4	67.7	0.76	0.71	85.1±0.2
CFGDistill samples	14.2	17.6	74.1	0.91	0.28	96.3±0.3
Canonical Samples	12.5	11.7	74.1	0.93	0.30	96.3±0.2

diffusion models, architectures, and samplers (Section F.2). Advances in either further improve scalability. Our goal is to demonstrate that core class semantics can be extracted from conditional diffusion models. We view improving efficiency as an important future work.

J QUANTITATIVE ANALYSIS OF CANONICAL SAMPLES

We provide quantitative metrics for a better understanding of the visual quality and representativeness of Canonical Samples generated by DiT (Peebles & Xie, 2023). We select 50000 images from the training set (Training samples), the samples used in CFGDistill (CFGDistill samples), and Canonical Samples, respectively. We compute the FID, sFID, Inception Score, Precision, and Recall, following the pipeline in (Dhariwal & Nichol, 2021). We report an additional metric called Classification Accuracy to show the discriminativeness of the samples. We use a pre-trained ResNet50, and randomly pick 10000 images from all sets of samples (Training, CFG, Canonical). We average the results over 3 seeds and report the mean and standard deviation as the error bars. The results are given in Table 3.

The reported values are different from random generative sampling in (Peebles & Xie, 2023), because we use DDIM inversion and decoding with a small number of time steps, instead of sampling from noise using a large number of function evaluations (NFEs). Note that our work does not focus on the visual quality of Canonical Samples, but instead on their semantics. Yet we observe that Canonical Samples have lower FID and sFID, indicating a better visual quality. The high precision and low recall suggest that samples that lie inside the data distribution while forming tight clusters, which is desired for class prototypes. We observe that the classifier gives higher classification accuracy on Canonical Samples, which means that they are more discriminative than the original samples.

We note that the pure CFG can also give similar results, particularly in Inception Score and Classification Accuracy. However, we emphasize that Canonical Samples and CFG yield different kinds of representative samples. Canonical Samples retain minimal class-irrelevant information, while CFG produces both class-irrelevant and class-related results. The visual results in Section Q support this claim. For example, in the "Tusker" class (Figure 44), pure CFG emphasizes "elephant" while Canonical Samples preserves both elephant and tusks, the joint defining attributes. The quantitative evidence is provided by our *CaDistill* experiments, in which CFGDistill yields models that suffer more from spurious features (Table 2) and generalize worse than the models produced by *CaDistill*.

K DETAILS OF THE PROOF-OF-CONCEPT EXPERIMENT

K.1 DATA GENERATION PROCESS

We adopt a hierarchical generative process described in Eq. 11 to generate 2D data points from two classes.

$$\begin{aligned}
 p(\mathbf{x}) &= p(y) p(\mathbf{x}_{\text{core}} | y) p(\mathbf{x}_{\text{var}} | \mathbf{x}_{\text{core}}), \\
 Y &\sim \text{Bernoulli}(\frac{1}{2}), \quad U \sim \mathcal{U}(-0.1, 0.1), \quad \varepsilon = (\varepsilon_x, \varepsilon_y)^\top \sim \mathcal{N}(\mathbf{0}, 0.01\mathbf{I}_2), \\
 \mathbf{s}(Y) &= \begin{cases} (0, 0)^\top, & Y = 0, \\ (4, 0)^\top, & Y = 1, \end{cases} \quad \mathbf{x}_{\text{core}} = (U, 0)^\top + \mathbf{s}(Y), \quad \mathbf{x}_{\text{var}} = \varepsilon, \\
 \mathbf{x} &= \mathbf{x}_{\text{core}} + \mathbf{x}_{\text{var}} + (3|\varepsilon_y|, 0)^\top.
 \end{aligned} \tag{11}$$

1998
 1999
 2000
 2001
 2002
 2003
 2004
 2005
 2006
 2007
 2008
 2009
 2010
 2011
 2012
 2013
 2014
 2015
 2016
 2017
 2018
 2019
 2020
 2021
 2022
 2023
 2024
 2025
 2026
 2027
 2028
 2029
 2030
 2031
 2032
 2033
 2034
 2035
 2036
 2037
 2038
 2039
 2040
 2041
 2042
 2043
 2044
 2045
 2046
 2047
 2048
 2049
 2050
 2051

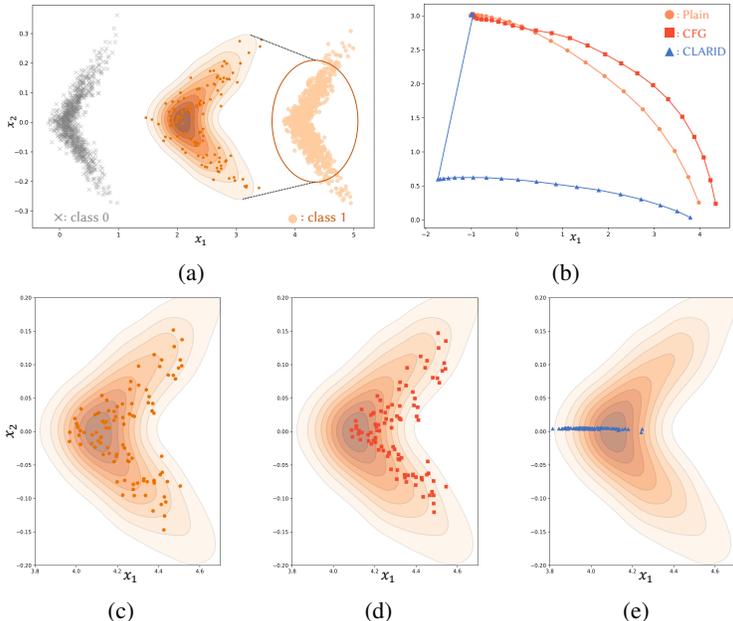


Figure 34: A toy example of CLARID. (a): The samples of class 0 and class 1; (b): Sampling trajectory of plain conditioning (Plain), classifier-free guidance (CFG), and CanoRep (CLARID) starting from $F_{inv}([4.0, 0.2])$. CanoRep (CLARID) orthogonalize the data latent code and class-irrelevant components encoded in the latent code in the CDM, yielding Canonical Samples; (c,d,e): The generated samples of Plain, CFG, and CLARID, respectively. CLARID produces Canonical Samples that lie on a 1D manifold inside class 1, offering an intuitive visual summary of the core class semantics.

Table 4: The hyperparameters used in our toy experiment in Section 3.3.

# data points	Epoch	Optimizer	Batch size	Learning rate	Weight decay	Label drop rate
1000	1000	Adam (Kingma & Ba, 2015)	128	1e-3	0	0.1

In our toy model, we have two classes. The process first generates class-specific samples on a segment $L = \{(x_1, 0) | 4y - 0.1 \leq x_1 \leq 4y + 0.1\}$, where $y \in \{0, 1\}$ denotes the class labels. It then adds class-independent noise to the points to generate the observed data. The samples from class 0 are included solely to introduce an inter-class contrast. We shift the x -axis value according to the y -axis variation to increase the complexity of the distribution.

K.2 ARCHITECTURE AND TRAINING DETAILS

We design a simple 3-layer multilayer-perceptron-based (MLP-based) diffusion model to model this 2D distribution. We set the hidden dimensionality to 80, and the dimensionality for both label and time step is 16. We use sinusoidal embeddings for uniquely encoding the 1000 time steps (Dosovitskiy et al., 2020). We train the diffusion model using the standard DDPM loss (Ho et al., 2020). The hyperparameters of training are given in Table 4.

K.3 CLARID AND BASELINE CONFIGURATION

When performing CLARID, we analytically compute the Jacobian of the network and calculate the corresponding singular vectors. We remove the first right singular vector at time step $t = 0.99T$. As a baseline, we take the latent codes obtained with F_{inv} and perform classifier-free guidance (CFG), steering the generation process toward regions with higher class-1 likelihood.

Table 5: The hyperparameters used in our experiments. We train a ResNet18 (He et al., 2016) on CIFAR10 and a ResNet50 on ImageNet. SGDM is SGD with momentum=0.9. *CaDistill* is effective with different training settings, as shown in Section M.2 and M.8.

\mathcal{D}	Epoch	Optimizer	Batch size	Learning rate	Weight decay	LR scheduler	LR decay rate
CIFAR10	200	SGDM	128	0.1	5e-4	Step 100,150	0.1
ImageNet100	100	SGDM	256	0.1	1e-4	Cosine	/
ImageNet	100	SGDM	512	0.1	1e-4	Cosine	/

K.4 RESULTS AND ANALYSIS

The results are shown in Figure 34. Notably, CLARID pushes most samples to a 1D manifold inside class 1, whereas CFG mainly steers the samples away from class 0. This low-dimensional manifold described by Canonical Samples can be regarded as a summarization of class 1 information in this case. The underlying structure revealed by Canonical Samples corresponds to one of the true generative processes for the observed data, which is the one used in our toy model. Reliably recovering the exact generative model is intractable due to the identifiability issue (Locatello et al., 2019). The solution generally requires extra inductive bias in modeling data distribution, which we leave for future work.

L TRAINING AND EVALUATION DETAILS OF *CaDistill*

It is often computationally demanding to compute the analytical Jacobian in Eq. 1. Because we only need the singular vectors, we adopt a well-established method, as in (Park et al., 2023a), to calculate only the top singular vectors of the Jacobian.

Regarding CaDistill, we perform all the experiments using the PyTorch platform. We provide the training hyperparameters in Table 5. The temperature parameter τ in *CaDistill* is fixed to 0.1 in all cases, as in (Khosla et al., 2020). On CIFAR10, we use random crop (`torchvision.transforms.RandomCrop(32, padding=4)`) and random horizontal flip (`torchvision.transforms.RandomHorizontalFlip`) for data augmentation. Note that the SupCon (Khosla et al., 2020) baseline uses a different data augmentation strategy and is trained for much longer epochs (1000). We follow the official code implementation (link) to reproduce this baseline on CIFAR10. On ImageNet and ImageNet100, we adopt the data augmentation used in Khosla et al. (2020), to improve the generalization performance of the trained model so that the baselines have meaningful results on the used generalization benchmarks. For example, the ResNet50 trained with random resized crop and random horizontal flip, as in He et al. (2016), will have 0% accuracy on ImageNet-A (Djolonga et al., 2021), while the baseline ResNet50 (Vanilla) in our experiments achieves 6.3%. The difference in data augmentation results in the performance difference on the ImageNet validation set between our baseline model and the one trained in He et al. (2016). However, we are not focusing on the clean performance in our settings. We use a single Nvidia A100 GPU for the CIFAR10 experiments, and four A100 GPUs for the ImageNet and ImageNet100 experiments. Training one model on CIFAR10 takes around 0.5 ~ 1 hour. Training on ImageNet takes around 30 hours for ResNet50, and 60 hours for ResNet152.

Our CaDistill involves the usage of CDM features, *i.e.*, Canonical Features. Instead of forwarding the CDM during training, which will lead to a large computational cost, we pre-compute the Canonical Features and load them during training, contributing to an efficient training pipeline (also see Section I). All Canonical Features are 1D vectors, which are the results of performing average pooling on the original feature map, as done in a previous work (Yang & Wang, 2023).

L.1 ADVERSARIAL ROBUSTNESS BENCHMARKING

We examine the adversarial robustness of different student models using four adversarial attacks, PGD (Madry et al., 2018), CW (Carlini & Wagner, 2017), APGD-DLR (Croce & Hein, 2020), and APGD-CE (Croce & Hein, 2020). The detailed settings are given in Table 6, 7, and 8 for CIFAR10, ImageNet, and ImageNet100 (AutoAttack), respectively. We choose PGD (Madry et al., 2018) since it is the most popular method for examining adversarial robustness (Xu et al., 2024). Moreover, we want to examine whether the drawback of cross-entropy loss can lead to false robustness (Croce &

Table 6: Hyperparameters in different adversarial attacks on CIFAR10.

	PGD	CW	APGD-DLR	APGD-CE
Max magnitude	2.0/255	/	2.0/255	2.0/255
Steps	5	5	5	5
Step size	0.5	/	/	/
κ	/	0.0	/	/
c	/	0.2	/	/
ρ	/	/	0.75	0.75
EOT	/	/	1	1

Table 7: Hyperparameters in different adversarial attacks on ImageNet.

	PGD	CW	APGD-DLR	APGD-CE
Max magnitude	0.33/255	/	0.33/255	0.33/255
Steps	5	5	5	5
Step size	0.5	/	/	/
κ	/	0.0	/	/
c	/	0.1	/	/
ρ	/	/	0.75	0.75
EOT	/	/	1	1

Hein, 2020). The step size in PGD can also largely affect the result. Hence, we choose the Auto-PGD family (Croce & Hein, 2020) to automatically decide the step size and incorporate the new Difference of Logits Ratio (DLR) loss, resulting in APGD-CE and APGD-DLR, respectively. We also want to include an optimization-based adversarial attack and thus select the CW (Carlini & Wagner, 2017) attack. The hyperparameters of different attacks are chosen to ensure a meaningful comparison between different models, avoiding the case in which all models have 0% accuracy after the attack. Our choices ensure a thorough test of the adversarial robustness in a white-box setting, revealing the multifacetedness of the adversarial robustness. Information about the generalization benchmarks is given in Section S.1.

L.2 GENERALIZATION ABILITY EVALUATION

On CIFAR10, we report the Top-1 accuracy on the CIFAR10-C (Hendrycks & Dietterich, 2018), CIFAR10.1 (Recht et al., 2018), CIFAR10.2 (Lu et al., 2020), as the metric for evaluating generalization. CIFAR10-C tests the out-of-distribution (OOD) generalization and CIFAR10.1, CIFAR10.2 evaluate the in-distribution (ID) one. On ImageNet, we show the Top-1 accuracy on ImageNet-C (IM-C) (Hendrycks & Dietterich, 2018), ImageNet-A (IM-A) (Djolonga et al., 2021), and ImageNet-Real (IM-Real) (Beyer et al., 2020). IM-C and IM-A are designed for benchmarking OOD generalization. IM-Real tests ID generalization.

Table 8: Hyperparameters of all attacks used in AutoAttack (Croce & Hein, 2020) in our ablation studies on ImageNet100.

	PGD	CW	APGD-DLR	APGD-CE
Max magnitude	0.5/255	/	0.5/255	0.5/255
Steps	5	10	5	5
Step size	0.5	/	/	/
κ	/	0.0	/	/
c	/	0.1	/	/
ρ	/	/	0.75	0.75
EOT	/	/	1	1

Table 9: Quantitative comparisons of *CaDistill* ImageNet (Deng et al., 2009) (ResNet-50). Adversarial robustness benchmarks: PGD (Madry et al., 2018), CW (Carlini & Wagner, 2017), APGD-DLR / APGD-CE (Croce & Hein, 2020); Evaluations of generalization ability: ImageNet-C (Hendrycks & Dietterich, 2018), ImageNet-A (Djolonga et al., 2021), ImageNet-Real (Beyer et al., 2020). Data_{DM} is the portion of data for which the DM acts as teacher. Higher is better. Values lower than the vanilla model are in red. See Section M for an analysis of the comparison between SimCLR and CLIP.

Model	Data_{DM}	Clean	PGD	CW	APGD-DLR	APGD-CE	IM-C	IM-A	IM-Real
Vanilla	/	75.9	15.6	13.7	17.2	16.7	45.9	6.3	82.8
SimCLR (2020b)	/	74.9	9.0	7.0	10.8	10.4	42.8	4.8	81.2
CLIP (2021)	/	68.2	7.0	8.8	0.2	1.1	21.2	11.4	74.4
DiffAug (2024)	100%	76.0	15.9	13.1	17.2	17.0	47.2	4.8	83.1
DMDistill	100%	75.7	15.7	14.1	17.0	16.7	43.6	5.0	82.8
CFGDistill	10%	75.7	20.8	20.3	20.8	21.4	45.6	6.0	82.7
<i>CaDistill</i>	10%	75.9	21.9	21.7	22.5	22.3	46.1	6.7	83.1

Due to the limits in computational resources, we do not report error bars in our experiments. During testing, we find that different runs of adversarial attacks result in similar performance. We test each adversarial attack 3 times with different seeds and find that the resulting performance has standard deviations all smaller than 0.05.

M MORE RESULTS ON IMAGENET

We provide additional results of *CaDistill* using ResNet50 on ImageNet in Table 9, 10. ResNet50 shows a similar trend as ResNet152 in the main paper. Hence, we perform extensive ablation studies and comparisons on ResNet50 in this section for efficiency.

Comparison to a vision-language model and a self-supervised learning method. In Table 9, 10, we compare models trained with *CaDistill* against two well-established baselines:

1. CLIP (Radford et al., 2021). It is a vision-language model that learns visual representations from image-text contrastive supervision. We use the ResNet50 vision backbone from the CLIP family. CLIP is pretrained on web-scale data rather than ImageNet. This difference in data distribution helps explain its strong performance on ImageNet-A, in which the images fall outside the ImageNet distribution but can be closer to CLIP’s pretraining distribution. For ImageNet1K probing, we use a nonlinear head (Linear-GELU-Linear) trained for 90 epochs with batch size 256, learning rate $1e-3$, AdamW, and no weight decay. The original paper reports 73.3% Top-1 using a single-batch L-BFGS linear probe, which is impractical in our setting due to memory requirements. Our linear head achieves 65.6% Top-1. Hence, we use the nonlinear head to improve the performance.
2. SimCLR (Chen et al., 2020a;b). It is a self-supervised method that learns invariances by contrasting two strongly augmented views of the same image. The augmentation is designed to be semantic-preserving. We use the official ResNet50 ($1\times$, no selective kernels) model (Chen et al., 2020b) trained on ImageNet1K and the official 100% finetuned checkpoint, which includes an ImageNet1K classification head.

Note that the data augmentation approaches and training epochs are not controlled in the comparison, because CLIP needs text as a source of supervision and SimCLR requires long training to be effective. Only the model architecture and size are controlled. Note that the CLIP ResNet50 has a slightly different architecture from the one used in our experiments and SimCLR. We emphasize that our approach complements mainstream supervised and self-supervised representation learning (Section 5 and O). These comparisons are provided for context rather than as a replacement claim: our goal is to offer a fresh perspective on learning representations, not to supplant existing paradigms.

M.1 THE PERFORMANCE OF \mathcal{L}_{dist} ALONE IN FEATURE DISTILLATION

In Section 4.1, we design a baseline experiment that distills the structure of the raw diffusion features into the representation space of the student network, which is based on \mathcal{L}_{dist} . Our design of \mathcal{L}_{dist}

2214 Table 10: Results on the Backgrounds Challenge (Xiao et al.) using ResNet50. Higher is better. See
 2215 Section M.5 for details.

2216 Model	2217 Original	2218 BG-Same	2219 BG-Rand	2220 Only-FG
2221 Vanilla	2222 96.0	2223 88.0	2224 81.1	2225 87.6
2226 SimCLR (2020b)	2227 94.9	2228 86.2	2229 79.7	2230 84.9
2231 CLIP (2021)	2232 87.7	2233 67.9	2234 56.2	2235 71.7
2236 DiffAug (2024)	2237 96.1	2238 87.5	2239 80.3	2240 87.4
2241 DMDistill	2242 96.3	2243 88.0	2244 80.6	2245 84.5
2246 CFGDistill	2247 96.3	2248 89.0	2249 82.6	2250 87.8
2251 CaDistill	2252 96.3	2253 89.0	2254 83.6	2255 88.5

2227 Table 11: Comparison between our CKA-based (Kornblith et al., 2019) \mathcal{L}_{dist} and other loss functions
 2228 in diffusion-based feature distillation.

2230 \mathcal{L}_{cano}	2231 Clean	2232 AutoAttack (2020)
2233 Vanilla (2016)	2234 86.5	2235 15.9
2236 FitNet (2023b; 2015; 2023)	2237 86.7	2238 16.4
2239 AT (2023b; 2023; 2017)	2240 86.6	2241 16.3
2242 RKD (2019; 2023)	2243 86.2	2244 16.4
2245 Ours	2246 87.3	2247 18.8

2238 differs from all previous works on diffusion-based feature distillation. We use a CKA (Kornblith et al.,
 2239 2019) metric for measuring the linear subspace alignment between the feature vectors of the student
 2240 and the teacher. CKA is invariant to isotropic scaling as well as orthonormal transformations. Such
 2241 an invariance lets us transfer the class-discriminative structure encoded in Canonical Features without
 2242 over-constraining the student’s own feature basis. Previous works focus on using three classical
 2243 feature distillation losses: (1) FitNet (Li et al., 2023b; Romero et al., 2015; Yang & Wang, 2023),
 2244 which is the L2 distance between student features and the teacher’s; (2) Attention transfer (AT) (Li
 2245 et al., 2023b; Yang & Wang, 2023; Zagoruyko & Komodakis, 2017), which distills the saliency
 2246 structure of the activation map to the student; (3) Relational knowledge distillation (RKD) (Park et al.,
 2247 2019; Yang & Wang, 2023), which aligns the relational representations of the samples between the
 2248 teacher and the student. However, in our case, these loss functions do not significantly contribute to
 2249 the student’s performance, as shown in Table 11. Note that in these experiments, the student features
 2250 and the teacher ones are one-to-one matching to mimic the typical feature distillation framework,
 2251 instead of the random strategy as we designed in **CaDistill**.

2252 Our CKA-based (Kornblith et al., 2019) feature distillation loss outperforms all previous designs
 2253 without introducing additional parameters during training. This is a novel loss function used in a
 2254 diffusion-based feature distillation framework, inspired by previous works (Dapello et al., 2023; Saha
 2255 et al., 2022; Zhou et al., 2024b). Performing knowledge transfer with the teacher and/or the student
 2256 being a ViT is still an open question (Yang et al., 2024; Yao et al., 2022), and can lead to a performance
 2257 drop in the student. Our \mathcal{L}_{dist} , however, achieves good performance in the diffusion-based settings.

2258 M.2 **CaDistill** IS EFFECTIVE WITH DIFFERENT DATA AUGMENTATION STRATEGIES

2259 We show that **CaDistill** is effective when the data augmentation strategy is different, demonstrating the
 2260 generalization of the paradigm. Specifically, the training lasts 120 epochs, and the data augmentations
 2261 are:

- 2262 • torchvision.transforms.RandomResizedCrop(224),
- 2263 • torchvision.transforms.RandomHorizontalFlip(),
- 2264 • torchvision.transforms.ColorJitter(0.3, 0.3, 0.3),

Table 12: Quantitative comparisons between *CaDistill* and baselines on ImageNet (Deng et al., 2009) with a ResNet50 (He et al., 2016), using a different training setting (Section M.2) from the one in Section L. Higher is better.

Model	Data _{DM}	Clean	PGD	CW	APGD-DLR	APGD-CE	IM-C	IM-A	IM-Real
Vanilla	/	76.6	17.3	13.5	18.8	17.7	40.6	3.4	83.3
CaDistill	10%	76.7	21.3	21.3	22.6	21.3	41.2	4.2	83.3

Table 13: Quantitative comparisons between *CaDistill*, and baselines on ImageNet (Deng et al., 2009) with a ResNet50 (He et al., 2016), on black-box adversarial robustness. Higher is better. Red is lower than the vanilla model. Data_{DM}: the portion of the subset on which the diffusion model serves as the teacher. DMDistill: Feature distillation by \mathcal{L}_{dist} on the whole dataset using a DiT model; CFGDistill: Using the framework of *CaDistill*, but replace CanoReps by samples with CFG from the CDM.

Model	Data _{DM}	Clean	Square (Andriushchenko et al., 2020)
Vanilla	/	75.9	23.5
DiffAug (Shama et al., 2024)	100%	76.0	20.7
DMDistill	100%	75.7	22.9
CFGDistill	10%	75.7	23.3
CaDistill	10%	75.9	25.6

The learning rate is 0.2, and the decay happens every 30 epochs with a decay rate of 0.1. The results are in Table 12. *CaDistill* yields a student that outperforms the vanilla model on all benchmarks, proving the effectiveness of our method in this case and implying its generalization capability.

M.3 *CaDistill* IMPROVES THE STUDENT’S BLACK-BOX ADVERSARIAL ROBUSTNESS

In Table 1, we demonstrate that *CaDistill* improves the student’s white-box robustness. Here, we show that *CaDistill* improves the student performance when facing black-box adversarial attacks. Specifically, we test all models on ImageNet using the Square attack (Andriushchenko et al., 2020), which is a black-box adversarial attack algorithm. We use L_{inf} metric with attacking budget 4.0/255.0, and the query number is 1000. To reduce computational costs, we randomly select 2000 samples from ImageNet to perform the evaluation. The result is given in Table 13.

M.4 EXTRACTING CLASS SEMANTICS USING CLASS TOKENS IN VISION TRANSFORMER

Our main claim on the application of CanoReps is that they distill the core class semantics of each category. Here, we investigate another mainstream approach to distill such information, which is the class token in vision transformer (ViT) (Dosovitskiy et al., 2020; Touvron et al., 2021; 2022). The class token performs the attention operation (Vaswani et al., 2017) to all spatial tokens, collecting the discriminative signals inside the feature map for classification. We adopt a challenging baseline network, DeiT-III-Huge (Touvron et al., 2022), which is a ViT model solely trained on ImageNet with the number of parameters (DeiT: 632.1M; DiT: 675M) and FLOPS (DeiT: 167.4G; DiT: 118.6G) matching the DiT (Peebles & Xie, 2023) used in our experiments. It is challenging because the DeiT-III-Huge model is trained with advanced data augmentation techniques and performs well on ImageNet classification tasks (Touvron et al., 2022), whereas DiT is trained with a plain horizontal flip augmentation and is not good at classification (Li et al., 2023a) (85.2 v.s. 77.5 Top-1 accuracy). We use our CKA-based \mathcal{L}_{dist} to align the representations of the student network to the class token in DeiT-III-Huge, termed DeiT_{dist}. All the settings are the same as used in Section 4.1 and L. The results are given in Table 14. Notably, the student trained with *CaDistill* outperforms the one trained with DeiT_{dist} in terms of clean accuracy and generalization. Achieving good performance in feature distillation between ViT and CNNs is still an open problem in the field (Yao et al., 2022). Despite this, our experiments control the architecture (both teachers are ViTs), the number of parameters, and the FLOPS. Moreover, the DeiT_{dist} can yield a student that outperforms DMDistill in terms of all adversarial attack benchmarks, which demonstrates the effectiveness of the method and validity of our

Table 14: Quantitative comparisons between *CaDistill*, and baselines on ImageNet (Deng et al., 2009) with a ResNet50 (He et al., 2016). Higher is better. Red is lower than the vanilla model. Data_{DM}: the portion of the subset on which the diffusion model serves as the teacher. DeiT_{dist}: Feature distillation by \mathcal{L}_{dist} on the whole dataset using the class token in an ImageNet-pretrained DeiT-III-Huge model (Touvron et al., 2022); DMDistill: Feature distillation by \mathcal{L}_{dist} on the whole dataset using a DiT model; CFGDistill: Using the framework of *CaDistill*, but replace CanoReps by samples with CFG from the CDM.

Model	Data _{DM}	Clean	PGD	CW	APGD-DLR	APGD-CE	IM-C	IM-A	IM-ReaL
Vanilla	/	75.9	15.6	13.7	17.2	16.7	45.9	6.3	82.8
DeiT _{dist} (Touvron et al., 2022)	100%	75.1	19.7	18.4	20.8	20.3	45.2	5.4	82.5
DMDistill	100%	75.7	15.7	14.1	17.0	16.7	43.6	5.0	82.8
CFGDistill	10%	75.7	20.8	20.3	20.8	21.4	45.6	6.0	82.7
<i>CaDistill</i>	10%	75.9	21.9	21.7	22.5	22.3	46.1	6.7	83.1

experiments. We believe investigating the difference between the mechanisms of how discriminative models and generative ones encode class information is an interesting future direction.

M.5 DETAILS OF THE BACKGROUND CHALLENGE

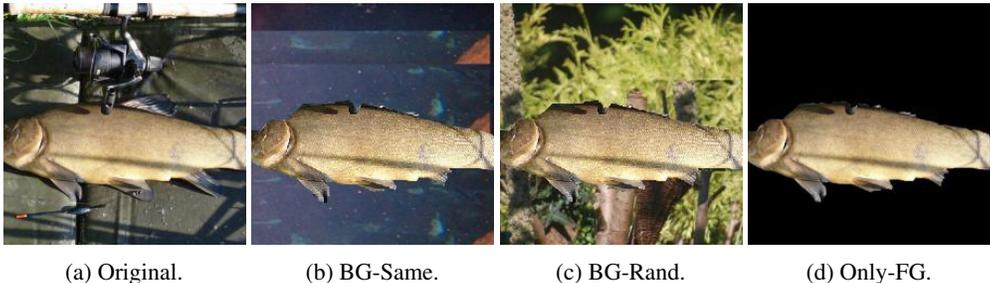


Figure 35: Samples from the Backgrounds Challenge (Xiao et al.). (a) Original: The original image. (b) BG-Same: Put a random background from the same class onto the image. (c) BG-Rand: Put a random background from a different class onto the image. (d) Only-FG: Discard the background and make it black.

In Section 4.1, we test the student model on the Backgrounds Challenge (Xiao et al.). Figure 35 illustrates its three variants: BG-Same re-uses a single background per class; BG-Rand pairs each foreground with randomly chosen backgrounds from other classes; Only-FG removes the background entirely. These operations preserve the foreground object while removing background cues. Hence, the performance in this test quantifies a model’s ability to rely on true class signals rather than spurious background correlations.

As shown in Table 2, *CaDistill* achieves the highest accuracy across all splits. CFGDistill matches *CaDistill* on the Original and BG-Same sets, but *CaDistill* outperforms it on BG-Rand and Only-FG. This gap indicates that the student trained with CFGDistill still uses background information shared within each class; when those backgrounds are shuffled or removed, its accuracy declines. The evidence suggests that CFG introduces label-correlated yet non-essential background signals into the training data, whereas *CaDistill* suppresses those signals and encourages the student to focus on the foreground object.

M.6 ON THE REPRODUCTION OF REPFUSION

RepFusion (Yang & Wang, 2023) proposes a novel framework for diffusion-based feature distillation. It uses a neural network for adaptively selecting the time step of feature extraction from the teacher DM. This neural network is trained using the REINFORCE (Williams, 1992) algorithm, using the task performance as the reward. The task performance is the classification accuracy. In this case, the neural network is non-linear and can directly decode the label conditioning in the CDM to maximize

2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429

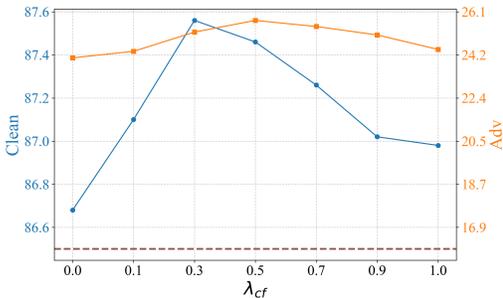


Figure 36: Ablation on λ_{cf} . An effective training requires a trade-off between \mathcal{L}_{align} and \mathcal{L}_{cano} , and necessitates both of them. We choose $\lambda_{cf} = 0.5$ to balance between the Clean accuracy and robustness. Brown is the baseline of both Clean and Adv.

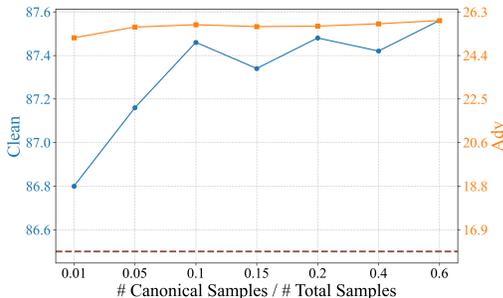


Figure 37: Ablation on the number of CanoReps. A small amount of CanoReps, e.g. 10%, is sufficient for achieving competitive performance. It implies the low-dimensionality property of the class manifolds inside CDMs, which is in line with previous findings (Wang et al.).

the reward, leading to a failed training. Hence, we reproduce the method on an unconditional DM. Despite this, we provide a strong baseline using CDM, DMDistill, which uses a feature distillation loss that can outperform all the loss functions used in RepFusion, as shown in Section M.1.

M.7 ABLATION STUDIES

We conduct ablation studies on ImageNet100, a 100-class subset of ImageNet. Previous studies (Cheng et al., 2023; Douillard et al., 2022; Xu et al., 2024; Yan et al., 2021; Yu et al., 2022) have shown that ImageNet100 serves as a representative subset of ImageNet1K. Hence, we can obtain representative results for the self-evaluation of the model while efficiently using our computational resources. All the ablations are based on a ResNet50 (He et al., 2016) model. Note that we still use an ImageNet-pretrained DiT (Peebles & Xie, 2023) as the teacher. We report the test accuracy on the ImageNet100 validation set as the clean accuracy (Clean), and the adversarial accuracy (Adv) under AutoAttack (Croce & Hein, 2020) that consists of PGD (Madry et al., 2018), CW (Carlini & Wagner, 2017), APGD-DLR (Croce & Hein, 2020), and APGD-CE (Croce & Hein, 2020).

M.7.1 NUMBER OF CANOREPS

Figure 37 shows the student performance when trained with different numbers of CanoReps. Remarkably, training with as little as 10% of the available CanoReps already yields near-optimal performance. The result suggests that a small data subset is enough to capture the core class semantics in CDMs, because those semantics lie on a low-dimensional manifold, which is consistent with earlier findings (Wang et al.).

M.7.2 THE NECESSITY OF \mathcal{L}_{align} AND \mathcal{L}_{cano} AND THEIR BALANCE

Our design includes on two complementary objectives: the alignment loss, \mathcal{L}_{align} , which pulls each sample towards CanoReps from its class, and the CanoRep separation loss, \mathcal{L}_{cano} , which drives the CanoReps of different classes apart. Figure 36 demonstrates the trade-off between the two. If \mathcal{L}_{align} is omitted ($\lambda_{cf} = 0$), samples remain distant from their canonical counterparts, preventing the student from learning the core semantics of each class. Conversely, dropping \mathcal{L}_{cano} ($\lambda_{cf} = 1$) can lead to CanoReps that collapse together, leaving the student unable to discriminate between categories. Therefore, the optimal choice requires a balance between them.

M.7.3 THE NECESSITY OF \mathcal{L}_{dist}

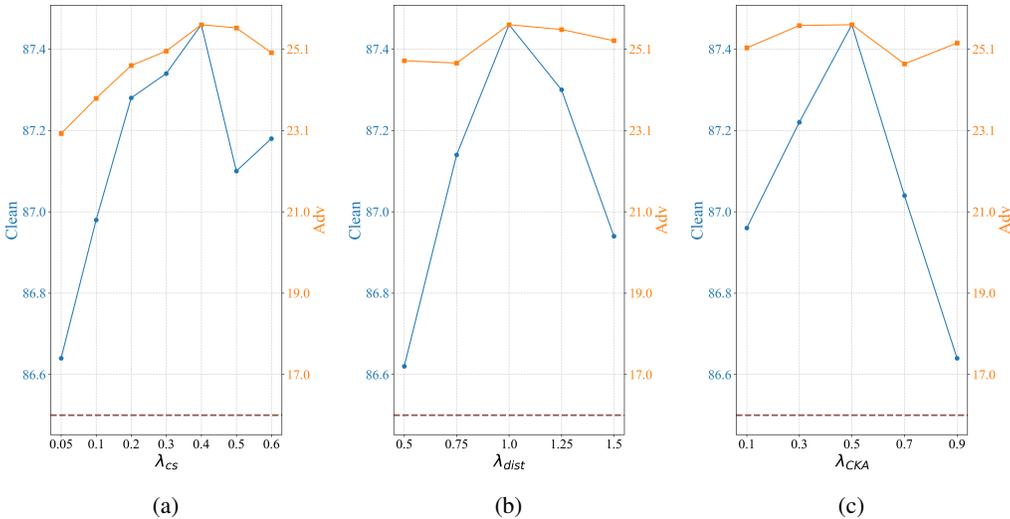
The CDM transfers the core class features using Canonical Features via \mathcal{L}_{dist} . Without this loss, the student can fail to learn the encoded features of the Canonical Samples, which can negatively affect the student’s clean accuracy and adversarial robustness, as shown in Table 15.

2430 Table 15: The ablation study of the effects of \mathcal{L}_{dist} on ImageNet (Deng et al., 2009) with a ResNet50
 2431 (He et al., 2016). Vanilla: The original student network. Data_{DM}: the portion of the subset on which
 2432 the diffusion model serves as the teacher. DMDistill: Feature distillation by \mathcal{L}_{dist} on the whole
 2433 dataset; CFGDistill: Using the framework of *CaDistill*, but replace Canonical Samples by samples
 2434 generated with CFG after F_{inv} , and use their corresponding features in the CDM. Higher is better.
 2435 **Green** is lower than the vanilla model. Without \mathcal{L}_{dist} , the student cannot learn the teacher’s encoding
 2436 of CanoReps, limiting the student’s adversarial robustness.

Model	Data _{DM}	Clean	PGD	CW	APGD-DLR	APGD-CE
Vanilla	/	75.9	15.6	13.7	17.2	16.7
DiffAug (Shama et al., 2024)	100%	76.0	15.9	13.1	17.2	17.0
DMDistill	100%	75.7	15.7	14.1	17.0	16.7
CFGDistill	10%	75.7	20.8	20.3	20.8	21.4
<i>CaDistill</i>	10%	75.9	21.9	21.7	22.5	22.3
No \mathcal{L}_{dist}	10%	75.6	20.3	19.3	20.5	21.9

2448 M.7.4 THE WEIGHTS OF LOSSES, $\lambda_{cs}, \lambda_{dist}, \lambda_{cka}$

2449 *CaDistill* involves 3 losses, $\lambda_{cs}, \lambda_{dist}, \lambda_{cka}$, each having its own weights. Here, we perform thorough
 2450 ablation studies on $\lambda_{cs}, \lambda_{dist}, \lambda_{cka}$ on ImageNet100. The results are given in Figure 38. We perform
 2451 the ablation study on one loss function by fixing the other weights to their own optimal values.
 2452 We empirically conclude this setting: $\lambda_{cs} = 0.4, \lambda_{dist} = 1.0, \lambda_{cka} = 0.5$, for all experiments on
 2453 ImageNet. For CIFAR10, we perform grid search over several parameter combinations and fix
 2454 $\lambda_{cs} = 0.2, \lambda_{dist} = 0.25, \lambda_{cka} = 0.5$.



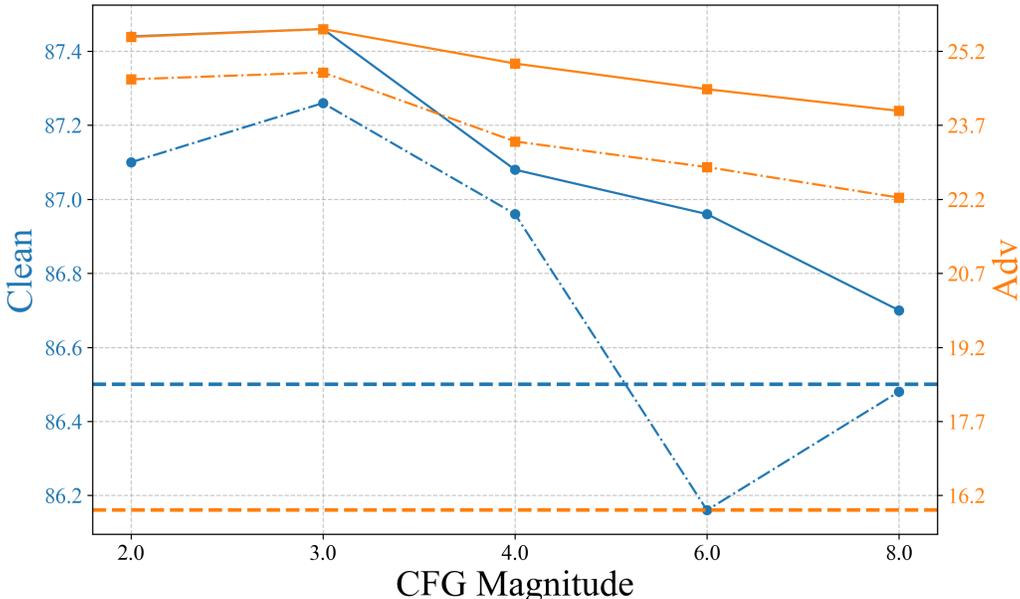
2474 Figure 38: Ablation studies on $\lambda_{cs}, \lambda_{dist}, \lambda_{cka}$, introduced in Section 4. We select $\lambda_{cs} = 0.4, \lambda_{dist} =$
 2475 $1.0, \lambda_{cka} = 0.5$ in all of our experiments on ImageNet.

2478 M.7.5 THE DESIGN OF \mathcal{L}_{cano}

2480 We design \mathcal{L}_{align} and \mathcal{L}_{cano} to both have push-together and pull-away effects, inspired by Khosla
 2481 et al. (2020). In Eq. equation 4, each Canonical Sample treats other Canonical Samples of
 2482 the same class (excluding itself) as positive examples. If a class happens to contribute only a single
 2483 Canonical Sample, no such positives exist. In that case, we optimize only the "pull-away" term—the
 denominator that separates the anchor from negatives in other classes, so the loss remains well-defined

2484 Table 16: The ablation study on using cross-entropy as \mathcal{L}_{cano} . Our design in Eq. 4 achieves better
 2485 results.

\mathcal{L}_{cano}	Clean	AutoAttack (Croce & Hein, 2020)
Vanilla (He et al., 2016)	86.5	15.9
Cross-entropy	86.8	25.3
Ours	87.5	25.7



2513 Figure 39: The ablation study on the magnitude of CFG used in our experiments on ImageNet. The
 2514 student is trained with **CaDistill** (solid lines) and CFGDistill (dash-dot lines). The dashed lines are the
 2515 baselines. Larger CFG magnitudes do not necessarily contribute to a better performance, indicating
 2516 that our design in **CaDistill** is not simply a converging prior on the features (Section G.3).

2517 and informative. In this case, \mathcal{L}_{cano} becomes:

$$2520 \mathcal{L}_{cano} = \frac{1}{b} \sum_{i=1}^b \log \sum_{k \neq i} \exp(\tilde{z}_i \cdot \tilde{z}_k / \tau). \tag{12}$$

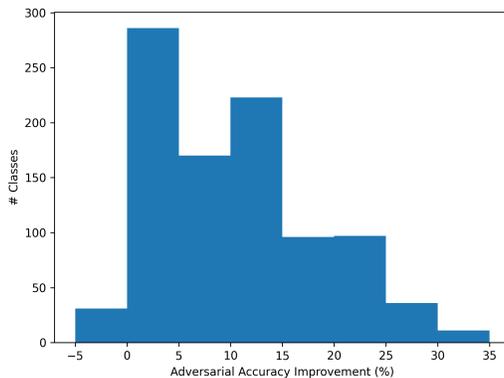
2524 We perform a simple ablation study on using cross-entropy for discriminating between CanoReps
 2525 from different classes, using ImageNet100. The results are given in Table 16. Our design achieves
 2526 better results in both clean accuracy and adversarial robustness, which is in line with the previous
 2527 claim (Khosla et al., 2020). In our case, CanoReps are far less than the original images and are
 2528 easier to classify, which can cause overfitting issues when using cross-entropy and lead to suboptimal
 2529 results.

2531 M.7.6 THE MAGNITUDE OF CLASSIFIER-FREE GUIDANCE

2533 On ImageNet, we use CFG after projecting away the extraneous directions. We perform an ablation
 2534 study on the CFG magnitude. The results are shown in Figure 39. Notably, larger CFG magnitudes
 2535 do not correspond to better performance. An overly large CFG scale can even worsen student
 2536 performance. This is because our **CaDistill** are not simply providing a converging prior over the
 2537 student features, as discussed in Section G.3. We choose the magnitude to be 3 for both **CaDistill** and
 CFGDistill.

2538 Table 17: Comparison between a vanilla Swin-Tiny (Liu et al., 2021) model, a FAN-Small (Zhou
 2539 et al., 2022) model, and the ones trained with *CaDistill*. Our method is effective with transformer
 2540 students. It also proves that *CaDistill* is effective with modern data augmentation techniques such as
 2541 Mixup (Zhang et al., 2018) and CutMix (Yun et al., 2019).

Model	Clean	AutoAttack (Croce & Hein, 2020)
Swin-Tiny (Liu et al., 2021)	81.8	9.3
<i>CaDistill</i>	84.4	13.8
FAN-Tiny (Zhou et al., 2022)	84.3	25.8
<i>CaDistill</i>	84.4	29.1



2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562 Figure 40: Class-wise adversarial robustness improvement of ResNet50 on ImageNet100 using
 2563 AutoAttack (Croce & Hein, 2020). *CaDistill* brings performance gain in almost all classes.

2564 2565 M.8 GENERALIZATION OF *CaDistill* TO DIFFERENT STUDENT ARCHITECTURES

2566 We demonstrate that *CaDistill* is effective when the student is a transformer architecture. Specifically,
 2567 we train a Swin-Tiny (Liu et al., 2021) model and a FAN-Tiny (Zhou et al., 2022) model on
 2568 ImageNet100. The result is given in Table 17. We train the networks using the same setting as
 2569 described in Section L, except that we follow the timm data augmentation with Mixup (Zhang et al.,
 2570 2018) and CutMix (Yun et al., 2019), and we have a 5-epoch learning rate warm-up.

2571 2572 M.9 CLASS-WISE ROBUSTNESS IMPROVEMENT

2573 We investigate whether the improvement is limited to a small set of classes or is spread across
 2574 the entire dataset. In Figure 40, we plot a histogram of the per-class percentage improvement in
 2575 adversarial accuracy under AutoAttack (Croce & Hein, 2020), using the same setting in Section M.7.
 2576 The histogram shows that almost all classes benefit from our method, indicating that the gains are
 2577 broadly distributed rather than concentrated in only a few categories.

2578 We further check whether the most-improved or least-improved classes align with specific super-
 2579 classes, such as "animals" or "artifacts". We list the 20 most improved classes:

2580
2581
2582 [`'hen-of-the-woods'`, `'racket'`, `'hip'`, `'go-kart'`, `'hyena'`, `'jacamar'`, `'Afghan hound'`, `'orangutan'`,
 2583 `'three-toed sloth'`, `'potter's wheel'`, `'lion'`, `'proboscis monkey'`, `'ostrich'`, `'steam locomotive'`, `'can-`
 2584 `'non'`, `'Komodo dragon'`, `'black and gold garden spider'`, `'partridge'`, `'gong'`, `'yurt'`];

2585 And 20 least improved classes:

2586
2587 [`'window screen'`, `'horned viper'`, `'beach wagon'`, `'rugby ball'`, `'hoopskirt'`, `'bib'`, `'Doberman'`,
 2588 `'kuvasz'`, `'Loafer'`, `'dial telephone'`, `'daisy'`, `'revolver'`, `'thunder snake'`, `'Kerry blue terrier'`, `'cocktail`
 2589 `'shaker'`, `'electric guitar'`, `'miniature pinscher'`, `'convertible'`, `'patio'`, `'Ibizan hound'`].

2590
2591 The results show that these classes are not clustered within only a small number of superclasses,
 suggesting that the robustness improvement is not biased toward any particular semantic group.

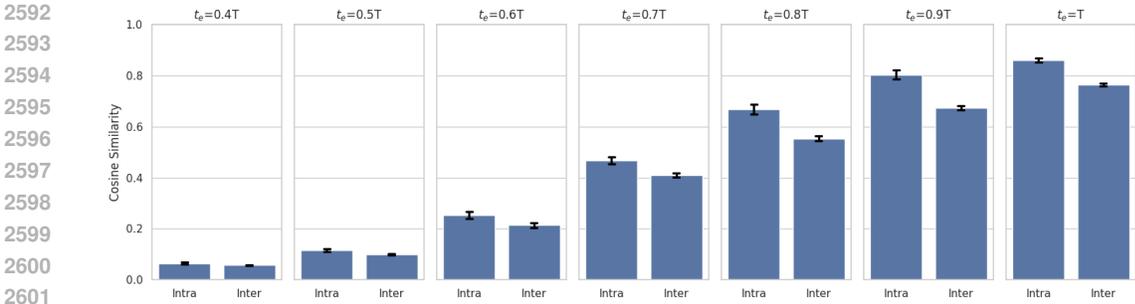


Figure 41: The cosine similarity between extraneous directions in both intra and inter class cases, with different t_e values. The error bar is the 95% CI.

M.10 THE VARIATION OF EXTRANEIOUS DIRECTIONS WITHIN AND BETWEEN CLASSES

We investigate the variation of extraneous directions within each class or across different classes. Specifically, we compute the cosine similarity between the top-10 extraneous directions that we obtain in the ImageNet20 experiment. In the intra-class case, we compute pair-wise cosine similarities between the top-10 directions and average them; In the inter-class case, we find an optimal one-to-one correspondence between the two sets of directions by solving the Hungarian matching problem, which maximizes the total pairwise similarity between matched directions. We conduct this experiment with $t_e \in \{0.4T, 0.5T, 0.6T, 0.7T, 0.8T, 0.9T, T\}$. The results are shown in Figure 41, with 95% confidence interval as the error bar. First, we observe that the similarity generally increases when t_e becomes larger. The values at $t_e = 0.9T$ and T are consistent with the prior work (Park et al., 2023b), validating our implementation. We also find that intra-class similarity is consistently higher than inter-class similarity. This is intuitive, as images from the same class often share similar backgrounds (e.g., most photos of sea fish are taken in the sea), leading to more similar extraneous directions within a class.

N SIGNIFICANCE OF THE QUANTITATIVE RESULTS AND DISCUSSION

N.1 SIGNIFICANCE OF THE IMPROVEMENTS BROUGHT BY *CaDistill*

Our proposed method, *CaDistill*, consistently improves the adversarial robustness and generalization ability of the student model. While the baseline methods can achieve better results on some benchmarks (e.g. the IM-C (Hendrycks & Dietterich, 2018) test on DiffAug (Shama et al., 2024)), they can worsen the student’s performance on other benchmarks (Red marks). This phenomenon reveals the established observation: the multifacetedness of robustness. Despite this, *CaDistill* still consistently improves the student’s performance, and it is the only one that is capable of doing so. This is neither a trivial nor marginal gain. Empirical and theoretical evidence suggest that there is a trade-off between different kinds of robustness (Moayeri et al., 2022; Rusak et al., 2020). Adversarial training can hurt corruption robustness (Rusak et al., 2020), while noise-based training weakens adversarial robustness (Table 1, row DiffAug). Adversarially robust models may even rely more on spurious cues (Moayeri et al., 2022). Our results and the CFGDistill rows exemplify these trade-offs, while *CaDistill* mitigates them by reducing reliance on spurious correlations.

N.2 ORIGIN OF IMPROVEMENTS

We argue that CFGDistill and *CaDistill* use different mechanisms to achieve adversarial robustness. First, CFGDistill shows a trade-off between adversarial robustness and generalization, a typical failure mode of current deep learning models (see above Significance of the improvements brought by *CaDistill*). In contrast, *CaDistill* overcomes such a trade-off. CFGDistill relies more on spurious correlations than *CaDistill*, as shown in Table 2. This shows that Canonical Samples matter as much as the loss. Our losses indeed contribute to adversarial robustness, but the structure of the involved samples plays an important role. We add a PlainDistill variant, where the samples are generated without CFG using DiT, to ablate sample effects. The results are given in Table 18.

Table 18: Quantitative comparisons between *CaDistill* and CFGDistill, PlainDistill on ImageNet (Deng et al., 2009) (ResNet-50). Adversarial robustness benchmarks: PGD (Madry et al., 2018), CW (Carlini & Wagner, 2017), APGD-DLR / APGD-CE (Croce & Hein, 2020); Evaluations of generalization ability : ImageNet-C (Hendrycks & Dietterich, 2018), ImageNet-A (Djolonga et al., 2021), ImageNet-ReaL (Beyer et al., 2020). Data_{DM} is the portion of data for which the DM acts as teacher. Higher is better. Values lower than the vanilla model are in red.

Model	Data_{DM}	Clean	PGD	CW	APGD-DLR	APGD-CE	IM-C	IM-A	IM-ReaL
Vanilla	/	75.9	15.6	13.7	17.2	16.7	45.9	6.3	82.8
PlainDistill	10%	75.4	17.3	15.7	18.3	18.7	45.3	5.8	82.6
CFGDistill	10%	75.7	20.8	20.3	20.8	21.4	45.6	6.0	82.7
<i>CaDistill</i>	10%	75.9	21.9	21.7	22.5	22.3	46.1	6.7	83.1

Adversarial robustness benefits from a converging prior in feature space (Pang et al.), that same-class samples are pulled together. Both CFGDistill and *CaDistill* provide such a prior, but they converge to different manifolds (Section 3.3): CFG’s manifold can encode non-causal structure, while Canonical samples encode the class core, yielding broader robustness. The ablation study in Section M.7.6 shows that excessive CFG magnitude degrades performance, demonstrating that CFG alone can impose the wrong structure (discussion in Section G.3; visual results in Figure 34). *CaDistill*’s advantage stems from both the Canonical Samples’ semantic structure and our loss design, producing comprehensive gains across various types of robustness. The results also suggest that Canonical Samples encode fundamentally different information, as the distillation from the two kinds of representations yields qualitatively (the trend of robustness improvement) and quantitatively (the absolute values of the robustness improvement) different models.

N.3 DATA EFFICIENCY OF *CaDistill*

The baseline methods that do not use our proposed feature distillation pipeline all necessitate access to the full dataset, while our methods achieve competitive performance with access only to 10% of the data. Given that the teacher model is often large and running it on the whole dataset can lead to high computational costs, this reduction in data dependency demonstrates the efficiency of our method.

N.4 DIFFERENCE OF THE RESULTS ON CIFAR10 AND IMAGENET

The performance trends differ in CIFAR10 and ImageNet (Table 1). On CIFAR10, even if simply distilling the raw diffusion features to the student via DMDistill can contribute to the performance on all benchmarks, while on ImageNet, all the baseline methods can perform worse than the vanilla model in certain cases. We assume that two factors can lead to such a difference. The first is that CIFAR10 is a simpler dataset than ImageNet. Most images in CIFAR10 contain purely foreground objects, while the images in ImageNet are much noisier and harder to classify. The evidence is that CIFAR10 is a nearly solved dataset, as the classification accuracy approaches 100% (Dosovitskiy et al., 2020), while the top models on ImageNet can achieve $\sim 90\%$ (Yu et al.).

More importantly, we identify a critical difference in the CDM on CIFAR10 and on ImageNet. That is, ImageNet-trained CDMs are typically trained in a low-resolution latent space of a pre-trained variational autoencoder (VAE) (Kingma & Welling, 2014; Rombach et al., 2022). This low-resolution space loses detailed information compared to the pixel space, reducing the discriminative power. To validate, we train a ResNet50 (He et al., 2016) in this latent space for image classification, receiving the input as the VAE-encoded images. Notably, the clean accuracy drastically drops from 86.5 to 76.6. We assume that the missing discriminative information can negatively affect the performance of all feature distillation methods, because the diffusion features lie in the low-resolution latent space. Due to the limits on computational resources, we leave the investigation on pixel-space DM on ImageNet as a future direction.

Table 19: The Pearson correlation r between the feature similarity matrices of different models and the ground truth class structure matrix (Huang et al., 2021) obtained by Wu-Palmer distance (Wu & Palmer, 1994) in Figure 18.

	r	p
ResNet50	0.19	< 0.0001
DeiT-Base-cls	0.05	0.6
DiT	0.31	< 0.0001
CFG	0.28	< 0.0001
Canonical Features	0.48	< 0.0001

O DISCUSSION AND POTENTIAL APPLICATIONS OF CLARID

CLARID identifies CanoReps that encode core class semantics while suppressing class-irrelevant information. Our *subtractive view contrasts and complements current representation learning research*, rely on supervised signals or contrastive objectives to learn semantics from inputs or to enforce invariances across inputs (e.g., vision–language models, VLMs, and self-supervised learning, SSL). We demonstrate that representations beneficial for robustness and generalization can be extracted from the full feature set used by the generator to synthesize the inputs. The full set itself, however, is less useful for high-level recognition as it encodes too much redundant information, as shown by our DMDistill experiments. This perspective does not contradict prior findings in diffusion-based representation learning showing gains in distillation (Li et al., 2023b; Yang & Wang, 2023), because they focus on low-level dense prediction tasks such as image segmentation and face landmark detection. In contrast, CLARID and *CaDistill* show that high-level semantics are also present in diffusion models. They simply need to be extracted by our principled method.

At a high level, this perspective aligns with generative similarity (Marjeh et al., 2024), where reliable object relationships and human-like class structure perception emerge from generative features, not purely from supervised or contrastive signals. The implication is that the generative process yields causal features that define the core semantics of objects (Section 3.3). We take a first practical step toward validating this claim and show promising gains in Table 1, 2. We further test whether our representations mirror human notions of class relationships. We evaluate the Pearson correlation r between the feature similarity matrix and the ground truth class structure matrix (Huang et al., 2021) in Figure 18, on ImageNet20. The results are given in Table 19. For ResNet50, we use the pre-trained model in torchvision and use the feature of the last layer after average pooling. DeiT-Base-cls is the class token of DeiT-Base (Touvron et al., 2021). These results indicate that Canonical Features capture inter-class relationships more faithfully than either CFG or supervised counterparts.

Diffusion inference produces reliable, human-like decisions (Jaini et al.; Li et al., 2023a) but is too compute-heavy for many deployments. *CaDistill* addresses this by transferring the CanoRep structure from the diffusion teacher to a lightweight student. This allows the student to have a similar semantic understanding. Acting as an amortized inference engine, the student approximates the teacher’s reasoning while avoiding its computational cost. This is important in resource-constrained settings like autonomous driving, where fast, reliable recognition is critical.

We also envision alternative practical use cases of CLARID on itself.

- **Semantic augmentation.** Canonical Samples capture each class’s core. By adding selected extraneous directions with controllable strength, we can generate high-diversity yet class-consistent variants. Such augmentations can complement underrepresented or long-tail categories during multi-modal pre-training. Figure 42 shows examples of the semantic augmentation.
- **Stronger visual grounding.** Canonical Samples suppress background clutter, offering clean contrasts for noisy real-world images. Combined with current visual pre-training pipelines, they can reduce reliance on spurious visual cues and improve grounding under domain shift.
- **Extract interpretable, compact class summaries and detect dataset bias.** Figure 6 and Section Q demonstrate that Canonical Samples capture the core class information, offering prototype-level class summaries. Canonical Samples highlight what the model considered irrelevant (background, co-occurring objects) versus indispensable to the class data, revealing potential biases of the dataset.

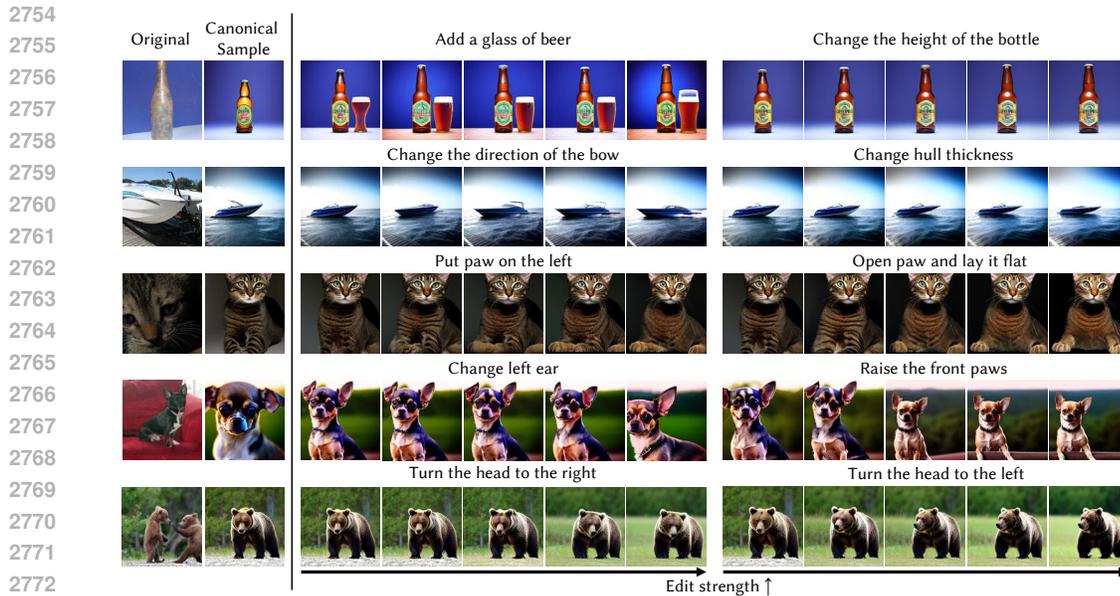


Figure 42: Semantic augmentation by adding extraneous directions back to CanoReps and decode, using Stable Diffusion 2.1. We show the effects of two extraneous directions for each sample, and different edit strengths. Note that all directions do not affect the class identity of the object.

For example, we spot that the Canonical Samples of the "Academic Gown" class (n02669723) always contain humans. After visual examination of the original class samples, we find that all images in this class have humans. This is an easily exploitable bias that can hinder model generalization. CLARID thus offers a diagnostic tool for dataset curation and auditing.

P ON MORE SOPHISTICATED FEATURE DISTILLATION FRAMEWORKS

In Section M.1, we demonstrate that our feature distillation loss outperforms the ones used in existing diffusion-based feature distillation frameworks. In this experiment, we use a single-layer distillation framework. That is, the alignment between the student and the teacher only happens at one layer, respectively. We do not consider more sophisticated feature distillation frameworks such as multi-layer alignments (Li et al., 2023b; Yang et al., 2024) due to limited computational resources. We believe investigating the combinations between *CaDistill* and different feature distillation frameworks is a promising future direction.

Q MORE VISUAL RESULTS

We provide more visualizations of CanoReps using Canonical Samples obtained from the DiT (Peebles & Xie, 2023) used in our *CaDistill* experiments on ImageNet (Deng et al., 2009), in Figure 43, 44, 45, 46, 47. All the classes are from ImageNet.

R BROADER IMPACT

CLARID introduces a new avenue into the field of DM research, focusing on interpreting the discriminative signals inside CDMs rather than directly probing the raw feature space. Our findings advance the interpretability of CDMs, contributing to safer usage of them. The application of CanoReps challenges the common assumption that the usage of DMs for discriminative tasks necessitates a large amount of data, providing new directions in diffusion-based feature distillation.



Figure 43: Visualizations of CanoReps using Canonical Samples obtained from the DiT (Peebles & Xie, 2023) used in our *CaDistill* experiments on ImageNet.



Figure 44: Visualizations of CanoReps using Canonical Samples obtained from the DiT (Peebles & Xie, 2023) used in our *CaDistill* experiments on ImageNet.

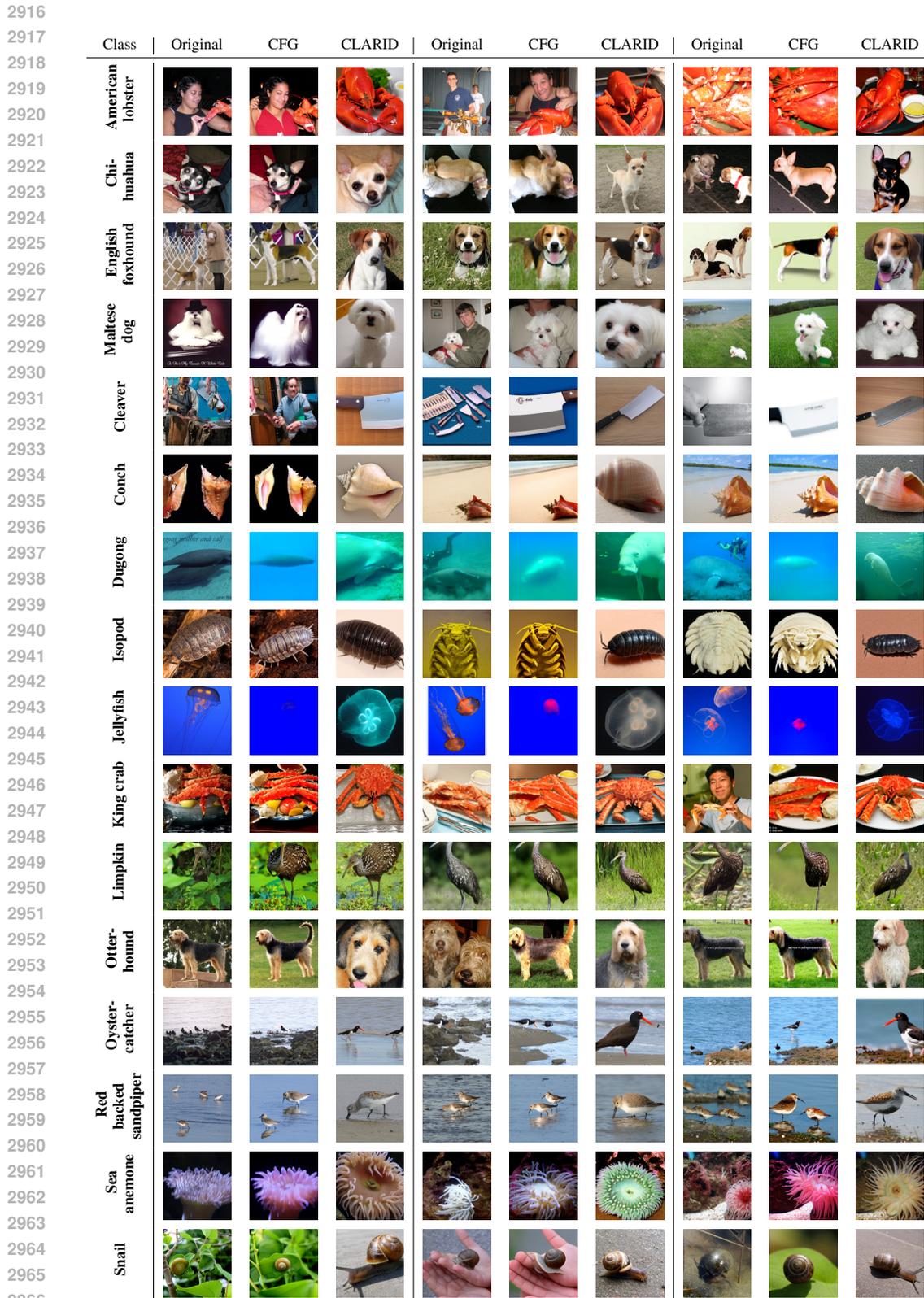


Figure 45: Visualizations of CanoReps using Canonical Samples obtained from the DiT (Peebles & Xie, 2023) used in our *CaDistill* experiments on ImageNet.

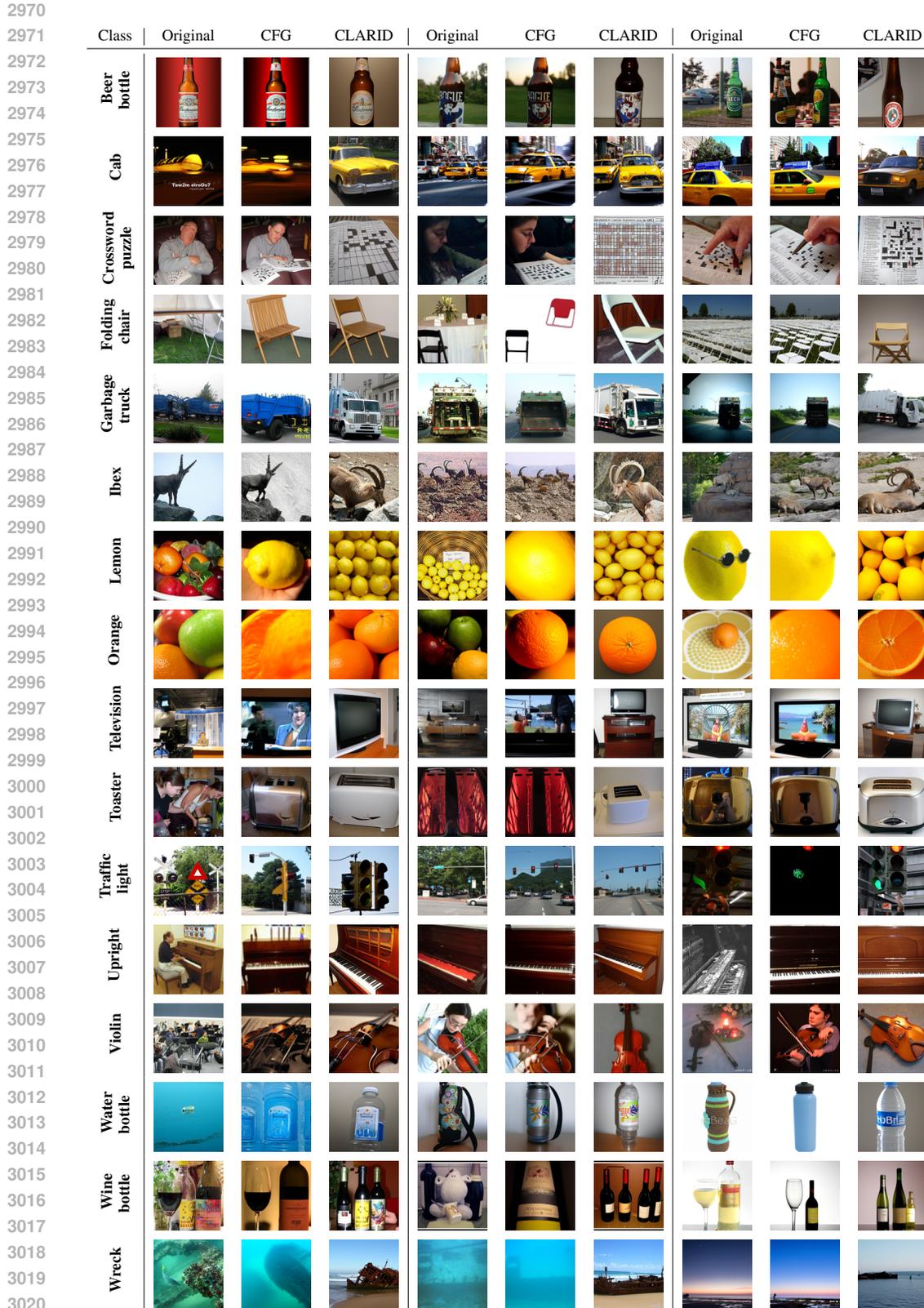


Figure 46: Visualizations of CanoReps using Canonical Samples obtained from the DiT (Peebles & Xie, 2023) used in our *CaDistill* experiments on ImageNet.

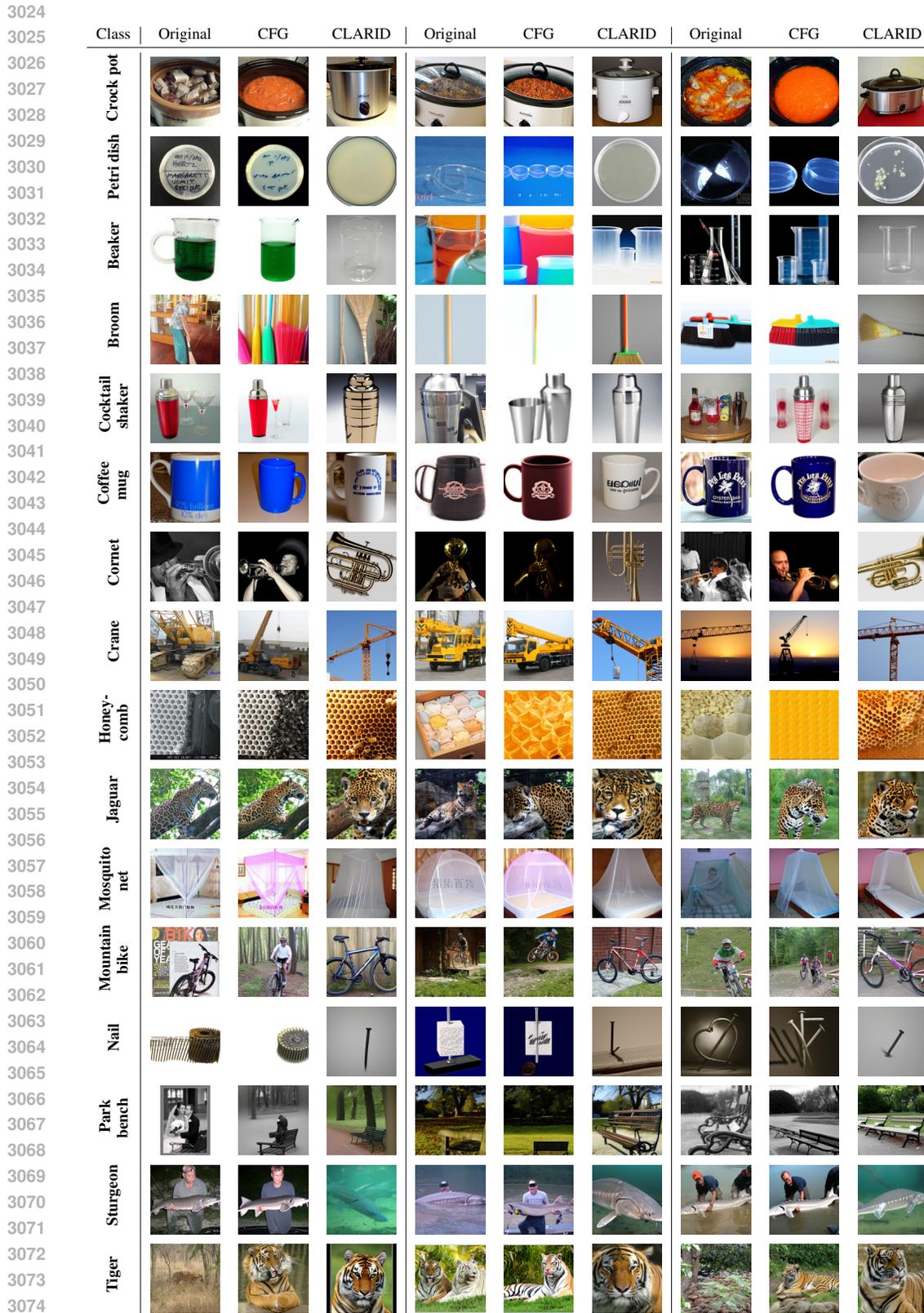


Figure 47: Visualizations of CanoReps using Canonical Samples obtained from the DiT (Peebles & Xie, 2023) used in our *CaDistill* experiments on ImageNet.

3078 S LICENSE INFORMATION
3079

3080 S.1 DATASETS INFORMATION AND LICENSE
3081

- 3082 • ImageNet1K (Deng et al., 2009). This dataset contains 1.28M training images and 50000 images
3083 for validation. We report the top1 accuracy on the 50000 validation images. License: Custom
3084 (research, non-commercial).
- 3085 • ImageNet-C (Hendrycks & Dietterich, 2018). This dataset contains 15 types of 2D image corruption
3086 types that are generated by different algorithms. Higher accuracy on this dataset indicates a more
3087 robust model against corrupted images. License: CC BY 4.0.
- 3088 • ImageNet-A (Djolonga et al., 2021). This dataset contains naturally existing adversarial examples
3089 that can drastically decrease the accuracy of ImageNet1K-trained CNNs. It is a 200-class subset of
3090 the ImageNet1K dataset. License: MIT license.
- 3091 • ImageNet Reassessed Labels (ImageNet-ReaL) (Beyer et al., 2020): This is a dataset with 50000
3092 reassessed labels of the ImageNet validation set, aiming at testing the in-distribution generalization
3093 ability of a classifier. License: Apache 2.0 License.
- 3094 • CIFAR10-C (Hendrycks & Dietterich, 2018). This dataset contains 15 types of 2D image corruption
3095 types that are generated by different algorithms. Higher accuracy on this dataset indicates a more
3096 robust model against corrupted images. License: CC BY 4.0.

3097
3098 S.2 MODEL AND CODE LICENSE
3099

- 3100 • Code for adversarial attacks (Kim, 2020): MIT License.
- 3101 • PyTorch Image Model (Wightman, 2019): Apache 2.0 License.
- 3102 • Diffusion Transformer (DiT) (Peebles & Xie, 2023): Attribution-NonCommercial 4.0 International.
- 3103 • Stable Diffusion 2.1 (Rombach et al., 2022): CreativeML Open RAIL++-M License.
- 3104 • EDM2 (Karras et al., 2024b): Creative Commons BY-NC-SA 4.0 license.
- 3105 • Supervised Contrastive Learning (Khosla et al., 2020): BSD 2-Clause License.
- 3106 • Swin (Liu et al., 2021): MIT License.
- 3107 • FAN (Zhou et al., 2022): Nvidia Source Code License-NC.
- 3108 • CIFAR10.1 (Recht et al., 2018): MIT License.

3109
3110
3111 T LARGE LANGUAGE MODELS USAGE

3112 We use Large Language Models to polish the writing.
3113
3114
3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131