# Deconstructing Denoising Diffusion Models for Self-Supervised Learning

Xinlei Chen[1]    Zhuang Liu[1]    Saining Xie[2]    Kaiming He[1]

[1]FAIR, Meta    [2]New York University

## Abstract

*In this study, we examine the representation learning abilities of Denoising Diffusion Models (DDM) that were originally purposed for image generation. Our philosophy is to deconstruct a DDM, gradually transforming it into a classical Denoising Autoencoder (DAE). This deconstructive procedure allows us to explore how various components of modern DDMs influence self-supervised representation learning. We observe that only a very few modern components are critical for learning good representations, while many others are nonessential. Our study ultimately arrives at an approach that is highly simplified and to a large extent resembles a classical DAE. We hope our study will rekindle interest in a family of classical methods within the realm of modern self-supervised learning.*

## 1. Introduction

Denoising is at the core in the current trend of generative models in computer vision and other areas. Popularly known as *Denoising Diffusion Models* (DDM) today, these methods [36, 37, 38, 23, 29, 11] learn a *Denoising Autoencoder* (DAE) [39] that removes noise of multiple levels driven by a diffusion process. These methods achieve impressive image generation quality, especially for high-resolution, photo-realistic images [33, 32]—in fact, these *generation* models are so good that they appear to have strong *recognition* representations for understanding the visual content.

While DAE is a powerhouse of today's generative models, it was originally proposed for learning representations [39] from data in a self-supervised manner. In today's community of representation learning, the arguably most successful variants of DAEs are based on "*masking noise*" [39], such as predicting missing text in languages (*e.g.*, BERT [10]) or missing patches in images (*e.g.*, MAE [21]). However, in concept, these masking-based variants remain significantly different from removing additive (*e.g.*, Gaussian) noise: while the masked tokens explicitly specify unknown *vs.* known content, no clean signal is available in the task of separating additive noise. Nevertheless, to-
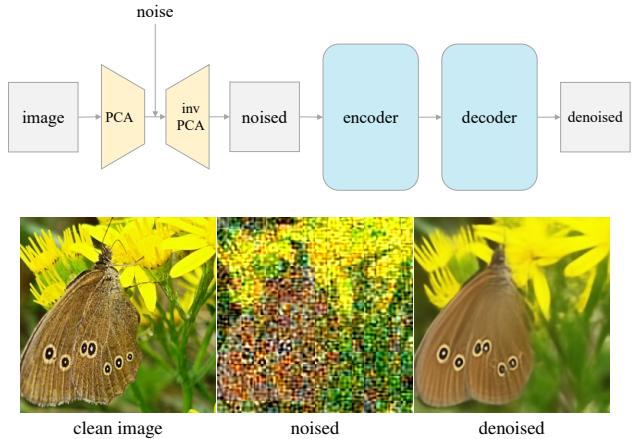


Figure 1. The latent Denoising Autoencoder (*l*-**DAE**) architecture we have ultimately reached, after a thorough exploration of deconstructing Denoising Diffusion Models (DDM) [23], with the goal of approaching the classical Denoising Autoencoder (DAE) [39] as much as possible. Here, the clean image (left) is projected onto a latent space using patch-wise PCA, in which noise is added (middle). It is then projected back to pixels via inverse PCA. An autoencoder is learned to predict a denoised image (right). This simple architecture largely resembles classical DAE (with the main difference that noise is added to the latent) and achieves competitive self-supervised learning performance.

day's DDMs for generation are dominantly based on additive noise, implying that they may learn representations without explicitly marking unknown/known content.

Most recently, there has been an increasing interest [40, 28] in inspecting the representation learning ability of DDMs. In particular, these studies directly take *off-the-shelf* pre-trained DDMs [23, 32, 11], which are originally purposed for generation, and evaluate their representation quality for recognition. They report encouraging results using these *generation-oriented* models. However, these pioneering studies obviously leave open questions: these off-the-shelf models were designed for generation, not recognition; it remains largely unclear whether the representation capability is gained by a denoising-driven process, or a diffusion-driven process.

In this work, we take a much deeper dive into the direction initialized by these recent explorations [40, 28]. Instead

1

of using an off-the-shelf DDM that is generation-oriented, we train models that are recognition-oriented. At the core of our philosophy is to *deconstruct* a DDM, changing it step-by-step into a classical DAE. Through this deconstructive research process, we examine every single aspect (that we can think of) of a modern DDM, with the goal of learning representations. This research process gains us new understandings on what are the critical components for a DAE to learn good representations.

Surprisingly, we discover that the main critical component is a tokenizer [33] that creates a *low-dimensional latent* space. Interestingly, this observation is largely *independent* of the specifics of the tokenizer—we explore a standard VAE [26], a patch-wise VAE, a patch-wise AE, and a patch-wise PCA encoder. We discover that it is the *low-dimensional* latent space, rather than the tokenizer specifics, that enables a DAE to achieve good representations.

Thanks to the effectiveness of PCA, our deconstructive trajectory ultimately reaches a simple architecture that is highly similar to the classical DAE (Fig. 1). We project the image onto a latent space using patch-wise PCA, add noise, and then project it back by inverse PCA. Then we train an autoencoder to predict a denoised image. We call this architecture "latent Denoising Autoencoder" (*l*-DAE).

Our deconstructive trajectory also reveals many other intriguing properties that lie between DDM and classical DAE. For one example, we discover that even using *a single noise level* (*i.e.*, not using the noise scheduling of DDM) can achieve a decent result with our *l*-DAE. The role of using multiple levels of noise is analogous to a form of data augmentation, which can be beneficial, but not an enabling factor. With this and other observations, we argue that the representation capability of DDM is mainly gained by the denoising-driven process, not a diffusion-driven process.

Finally, we compare our results with previous baselines. On one hand, our results are substantially better than the off-the-shelf counterparts (following the spirit of [40, 28]): this is as expected, because these are our starting point of deconstruction. On the other hand, our results fall short of baseline contrastive learning methods (*e.g.*, [7]) and masking-based methods (*e.g.*, [21]), but the gap is reduced. Our study suggests more room for further research along the direction of DAE and DDM.

## 2. Related Work

In the history of machine learning and computer vision, the generation of images (or other content) has been closely intertwined with the development of unsupervised or self-supervised learning. Approaches in generation are conceptually forms of un-/self-supervised learning, where models were trained without labeled data, learning to capture the underlying distributions of the input data.

There has been a prevailing belief that the ability of

a model to generate high-fidelity data is indicative of its potential for learning good representations. Generative Adversarial Networks (GAN) [18], for example, have ignited broad interest in adversarial representation learning [13, 12]. Variational Autoencoders (VAEs) [26], originally conceptualized as generative models for approximating data distributions, have evolved to become a standard in learning localized representations ("tokens"), *e.g.*, VQVAE [30] and variants [16]. Image inpainting [2], essentially a form of conditional image generation, has led to a family of modern representation learning methods, including Context Encoder [31] and Masked Autoencoder (MAE) [21].

Analogously, the outstanding generative performance of Denoising Diffusion Models (DDM) [36, 37, 38, 23, 11] has drawn attention for their potential in representation learning. Pioneering studies [40, 28] have begun to investigate this direction by evaluating existing pre-trained DDM networks. However, we note that while a model's generation capability suggests a certain level of understanding, it does not necessarily translate to representations useful for downstream tasks. Our study delves deeper into these issues.

On the other hand, although Denoising Autoencoders (DAE) [39] have laid the groundwork for autoencoding-based representation learning, their success has been mainly confined to scenarios involving masking-based corruption (*e.g.*, [21, 41, 17, 5]). To the best of our knowledge, little or no recent research has reported results on classical DAE variants with additive Gaussian noise, and we believe that the underlying reason is that a simple DAE baseline (Fig. 2(a)) performs poorly[1] (*e.g.*, in ~20% Fig. 5).

## 3. Background: Denoising Diffusion Models

Our deconstructive research starts with a Denoising Diffusion Model (DDM) [36, 37, 38, 23, 11]. We briefly describe the DDM we use, following [11, 32].

A diffusion process starts from a clean data point ($z_0$) and sequentially adds noise to it. At a specified time step $t$, the noised data $z_t$ is given by:

$$z_t = \gamma_t z_0 + \sigma_t \epsilon \qquad (1)$$

where $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ is a noise map sampled from a Gaussian distribution, and $\gamma_t$ and $\sigma_t$ define the scaling factors of the signal and of the noise, respectively. By default, it is set $\gamma_t^2 + \sigma_t^2 = 1$ [29, 11].

A denoising diffusion model is learned to remove the noise, conditioned on the time step $t$. Unlike the original DAE [39] that predicts a clean input, the modern DDM [23, 29] often predicts the noise $\epsilon$. Specifically, a

---

[1]According to the authors of MoCo [20] and MAE [21], significant effort has been devoted to DAE baselines during the development of those works, following the best practice established. However, it has not led to meaningful results (<20% accuracy).
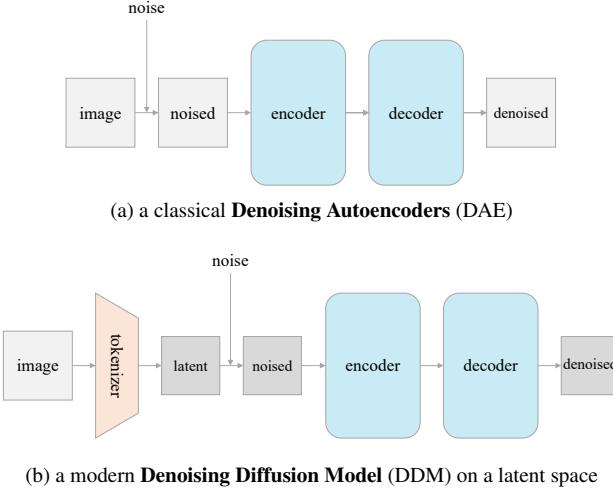
(a) a classical **Denoising Autoencoders** (DAE)



(b) a modern **Denoising Diffusion Model** (DDM) on a latent space

Figure 2. **A classical DAE and a modern DDM**. (a) A classical DAE that adds and predicts noise on the image space. (b) State-of-the-art DDMs (*e.g.*, LDM [33], DIT [32]) that operate on a latent space, where the noise is added and predicted.

loss function in this form is minimized:

$$\|\epsilon - \mathtt{net}(z_t)\|^2 \qquad (2)$$

where $\mathtt{net}(z_t)$ is the network output. The network is trained for multiple noise levels given a noise *schedule*, conditioned on the time step $t$. In the generation process, a trained model is iteratively applied until it reaches the clean signal $z_0$.

DDMs can operate on two types of input spaces. One is the original pixel space [11], where the raw image $x_0$ is directly used as $z_0$. The other option is to build DDMs on a *latent* space produced by a *tokenizer*, following [33]. See Fig. 2(b). In this case, a pre-trained tokenizer $f$ (which is often another autoencoder, *e.g.*, VQVAE [30]) is used to map the image $x_0$ into its latent $z_0 = f(x_0)$.

**Diffusion Transformer (DiT).** Our study begins with the Diffusion Transformer (DiT) [32]. We choose this Transformer-based DDM for several reasons: (i) Unlike other UNet-based DDMs [11, 33], Transformer-based architectures can provide fairer comparisons with other self-supervised learning baselines driven by Transformers (*e.g.*, [7, 21]); (ii) DiT has a clearer distinction between the encoder and decoder, while a UNet's encoder and decoder are connected by skip connections and may require extra effort on network surgery when evaluating the encoder; (iii) DiT trains much faster than other UNet-based DDMs (see [32]) while achieving better generation quality.

We use the DiT-Large (**DiT-L**) variant [32] as our DDM baseline. In DiT-L, the encoder and decoder *put together* have the size of ViT-L [15] (24 blocks). We evaluate the representation quality (linear probe accuracy) of the *encoder*, which has 12 blocks, referred to as "$\frac{1}{2}$L" (half large).

**Tokenizer.** DiT instantiated in [32] is a form of Latent Diffusion Models (LDM) [33], which uses a VQGAN tokenizer [16]. Specifically, this VQGAN tokenizer transforms the $256 \times 256 \times 3$ input image (height×width×channels) into a $32 \times 32 \times 4$ latent map, with a stride of 8.

**Starting baseline.** By default, we train the models for 400 epochs on ImageNet [9] with a resolution of $256 \times 256$ pixels. Implementation details are in Sec. A.

Our DiT baseline results are reported in Tab. 1 (line 1). With DiT-L, we report a linear probe accuracy of **57.5%** using its $\frac{1}{2}$L encoder. The generation quality (Fréchet Inception Distance [22], FID-50K) of this DiT-L model is **11.6**. This is the starting point of our destructive trajectory.

Despite differences in implementation details, our starting point conceptually follows recent studies [40, 28] (more specifically, DDAE [40]), which evaluate off-the-shelf DDMs under the linear probing protocol.

## 4. Deconstructing Denoising Diffusion Models

Our deconstruction trajectory is divided into three stages. We first adapt the generation-focused settings in DiT to be more oriented toward self-supervised learning (Sec. 4.1). Next, we deconstruct and simplify the tokenizer step by step (Sec. 4.2). Finally, we attempt to reverse as many DDM-motivated designs as possible, pushing the models towards a classical DAE [39] (Sec. 4.3). We summarize our learnings from this deconstructing process in Sec. 4.4.

### 4.1. Reorienting DDM for Self-supervised Learning

While a DDM is conceptually a form of a DAE, it was originally developed for the purpose of image generation. Many designs in a DDM are oriented toward the generation task. Some designs are *not legitimate* for self-supervised learning (*e.g.*, class labels are involved); some others are not necessary if visual quality is not concerned. In this subsection, we reorient our DDM baseline for the purpose of self-supervised learning, summarized in Tab. 1.

**Remove class-conditioning.** A high-quality DDM is often trained with conditioning on *class labels*, which can largely improve the generation quality. But the usage of class labels is simply not legitimate in the context of our self-supervised learning study. As the first step, we remove class-conditioning in our baseline.

Surprisingly, removing class-conditioning substantially improves the linear probe accuracy from 57.5% to 62.1% (Tab. 1), even though the generation quality is greatly hurt as expected (FID from 11.6 to 34.2). We hypothesize that directly conditioning the model on class labels may reduce the model's demands on encoding the information related to class labels. Removing the class-conditioning can force the model to learn more semantics.

3

|                              | acc. (↑) | FID (↓) |
|------------------------------|----------|---------|
| DiT baseline                 | 57.5     | 11.6    |
| + remove class-conditioning  | 62.5     | 30.9    |
| + remove VQGAN perceptual loss | 58.4   | 54.3    |
| + remove VQGAN adversarial loss | 59.0  | 75.6    |
| + replace noise schedule     | 63.4     | 93.2    |

Table 1. **Reorienting DDM for self-supervised learning**. We begin with the DiT [32] baseline and evaluate its linear probe accuracy (acc.) on ImageNet. Each line is based on a modification of the immediately preceding line. The entries in gray, in which class labels are used, are not legitimate results for self-supervised learning. See Sec. 4.1 for description.

**Deconstruct VQGAN.** In our baseline, the VQGAN tokenizer, presented by LDM [33] and inherited by DiT, is trained with multiple loss terms: (i) autoencoding reconstruction loss; (ii) KL-divergence regularization loss [33];[2] (iii) perceptual loss [44] based on a *supervised* VGG net [35] trained for ImageNet classification; and (iv) adversarial loss [18, 16] with a discriminator. We ablate the latter two terms in Tab. 1.

As the **perceptual loss** [44] involves a supervised pretrained network, using the VQGAN trained with this loss is not legitimate. Instead, we train another VQGAN tokenizer [33] in which we remove the perceptual loss. Using this tokenizer *reduces* the linear probe accuracy significantly from 62.5% to 58.4% (Tab. 1), which, however, provides the first legitimate entry thus far. This comparison reveals that *a tokenizer trained with the perceptual loss (with class labels) in itself provides semantic representations*. We note that the perceptual loss is not used from now on, in the remaining part of this paper.

We train the next VQGAN tokenizer that further removes the **adversarial loss**. It slightly increases the linear probe accuracy from 58.4% to 59.0% (Tab. 1). With this, our tokenizer at this point is essentially a VAE, which we move on to deconstruct in the next subsection. We also note that removing either loss harms generation quality.

**Replace noise schedule.** In the task of generation, the goal is to progressively turn a noise map into an image. As a result, the original noise schedule spends many time steps on very noisy images (Fig. 3). This is not necessary if our model is not generation-oriented.

We study a simpler noise schedule for the purpose of self-supervised learning. Specifically, we let $\gamma_t^2$ linearly decay in the range of $1 > \gamma_t^2 \geq 0$ (Fig. 3). This allows the model to spend more capacity on cleaner images. This change greatly improves the linear probe accuracy from 59.0% to 63.4% (Tab. 1), suggesting that the original schedule focuses too much on noisier regimes. On the other hand, as expected, doing so further hurts the generation ability, leading to a FID of 93.2.

---

[2]The KL form in [33] does not perform explicit vector quantization (VQ), interpreted as "the quantization layer absorbed by the decoder" [33].



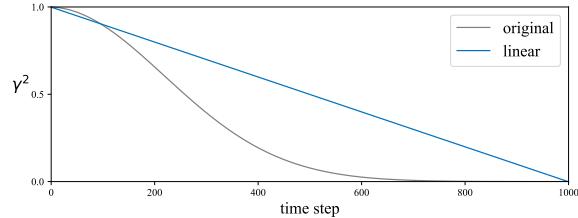Figure 3. **Noise schedules**. The original schedule [23, 32], which sets $\gamma_t^2 = \Pi_{s=1}^t (1 - \beta_s)$ with a linear schedule of $\beta$, spends many time steps on very noisy images (small $\gamma$). Instead, we use a simple schedule that is linear on $\gamma^2$, which provides less noisy images.

**Summary.** Overall, the results in Tab. 1 reveal that *self-supervised learning performance is not correlated to generation quality*. The representation capability of a DDM is not necessarily the outcome of its generation capability.

### 4.2. Deconstructing the Tokenizer

Next, we further deconstruct the VAE tokenizer by making substantial simplifications. We compare the following four variants of autoencoders as the tokenizers, each of which is a simplified version of the preceding one:

- **Convolutional VAE.** Our deconstruction thus far leads us to a VAE tokenizer. As common practice [26, 33], the encoder $f(\cdot)$ and decoder $g(\cdot)$ of this VAE are deep convolutional (conv) neural networks [27]. This convolutional VAE minimizes the following loss function:

$$\|x - g(f(x))\|^2 + \mathbb{KL}\left[f(x)|\mathcal{N}\right].$$

Here, $x$ is the input image of the VAE. The first term is the reconstruction loss, and the second term is the Kullback-Leibler divergence [3, 16] between the latent distribution of $f(x)$ and a unit Gaussian distribution.

- **Patch-wise VAE.** Next we consider a simplified case in which the VAE encoder and decoder are both *linear* projections, and the VAE input $x$ is a *patch*. The training process of this *patch-wise VAE* minimizes this loss:

$$\|x - U^T V x\|^2 + \mathbb{KL}\left[V x|\mathcal{N}\right].$$

Here $x$ denotes a patch flattened into a $D$-dimensional vector. Both $U$ and $V$ are $d \times D$ matrixes, where $d$ is the dimension of the latent space. We set the patch size as $16 \times 16$ pixels, following [15].

- **Patch-wise AE.** We make further simplification on VAE by removing the regularization term:

$$\|x - U^T V x\|^2.$$

As such, this tokenizer is essentially an autoencoder (AE) on patches, with the encoder and decoder both being linear projections.
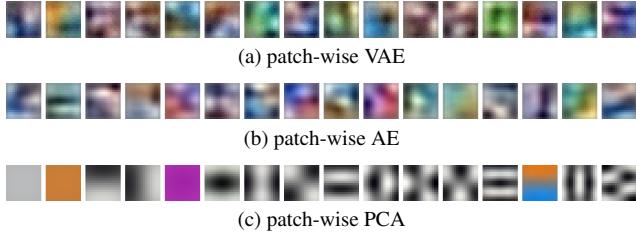
4

(a) patch-wise VAE

(b) patch-wise AE

(c) patch-wise PCA

Figure 4. **Visualization of the patch-wise tokenizer**. Each filter corresponds to a row of the linear projection matrix $V$ ($d \times D$), reshaped to $16 \times 16 \times 3$ for visualization. Here $d$=16.

• **Patch-wise PCA.** Finally, we consider a simpler variant which performs Principal Component Analysis (PCA) on the patch space. It is easy to show that PCA is equivalent to a special case of AE:

$$\|x - V^T V x\|^2.$$

in which $V$ satisfies $VV^T=I$ ($d \times d$ identity matrix). The PCA bases can be simply computed by eigen-decomposition on a large set of randomly sampled patches, requiring no gradient-based training.

Thanks to the simplicity of using patches, for the three patch-wise tokenizers, we can visualize their filters in the patch space (Fig. 4).

Tab. 2 summarizes the linear probe accuracy of DiT using these four variants of tokenizers. We show the results w.r.t. the *latent dimension "per token"*.[3] We draw the following observations.
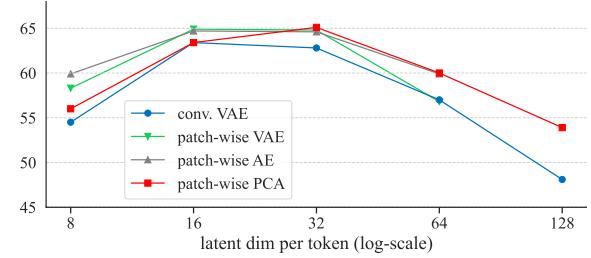
***Latent dimension of the tokenizer is crucial for DDM to work well in self-supervised learning.***

As shown in Tab. 2, all four variants of tokenizers exhibit similar trends, despite their differences in architectures and loss functions. Interestingly, the optimal dimension is relatively low ($d$ is 16 or 32), even though the full dimension per patch is much higher ($16 \times 16 \times 3$=768).

Surprisingly, the convolutional VAE tokenizer is neither necessary nor favorable; instead, all patch-based tokenizers, in which each patch is encoded *independently*, perform similarly with each other and consistently outperform the Conv VAE variant. In addition, the KL regularization term is *unnecessary*, as both the AE and PCA variants work well.

To our further surprise, *even the PCA tokenizer works well*. Unlike the VAE or AE counterparts, the PCA tokenizer does *not* require gradient-based training. With pre-computed PCA bases, the application of the PCA tokenizer

---

[3]For patch-wise VAE/AE/PCA (patch stride is 16), we treat each patch as a token, so the latent dimension is simply $d$ for each patch. For the default convolutional VAE that has a stride of 8, the DiT implementation [32] treats each $2 \times 2$ patch on the latent space as a "token"; as a result, its latent dimension "per token" should be multiplied by 4 for calibration.



| latent dim. $d$ | 8 | 16 | 32 | 64 | 128 |
|---|---|---|---|---|---|
| conv. VAE (baseline) | 54.5 | **63.4** | 62.8 | 57.0 | 48.1 |
| patch-wise VAE | 58.3 | **64.9** | 64.8 | 56.8 | - |
| patch-wise AE | 59.9 | **64.7** | 64.6 | 59.9 | - |
| patch-wise PCA | 56.0 | 63.4 | **65.1** | 60.0 | 53.9 |

Table 2. **Linear probe accuracy *vs*. latent dimension**. With a DiT model, we study four variants of tokenizers for computing the latent space. We vary the dimensionality $d$ (*per token*) of the latent space. The table is visualized by the plot above. **All four variants of tokenizers exhibit similar trends**, despite their differences in architectures and loss functions. The 63.4% entry of "conv. VAE" is the same entry as the last line in Tab. 1.
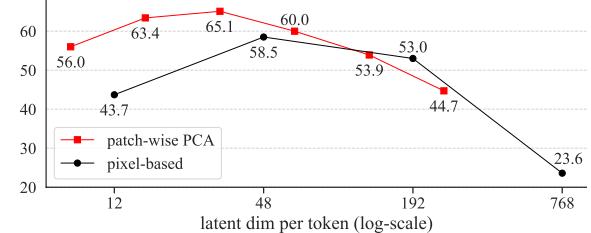


Figure 5. Linear probe results of the **pixel-based tokenizer**, operated on an image size of 256, 128, 64, and 32, respectively with a patch size of 16, 8, 4, 2. The "latent" dimensions of these tokenized spaces are 768, 192, 48, and 12 per token. Similar to other tokenizers we study, this pixel-based tokenizer exhibits a similar trend: a relatively small dimension of the latent space is optimal.

is analogous to a form of image pre-processing, rather than a "network architecture". The effectiveness of a PCA tokenizer largely helps us push the modern DDM towards a classical DAE, as we will show in the next subsection.

***High-resolution, pixel-based DDMs are inferior for self-supervised learning.***

Before we move on, we report an extra ablation that is consistent with the aforementioned observation. Specifically, we consider a "naïve tokenizer" that performs *identity mapping* on patches extracted from resized images. In this case, a "token" is the flatten vector consisting all *pixels* of a patch.

In Fig. 5, we show the results of this "pixel-based" tokenizer, operated on an image size of 256, 128, 64, and 32, respectively with a patch size of 16, 8, 4, 2. The "latent" dimensions of these tokenized spaces are 768, 192, 48, and 12 per token. In all case, the sequence length of the Transformer is kept unchanged (256).

| | acc. |
|---|---|
| patch-wise PCA baseline | 65.1 |
| + predict clean data (rather than noise) | 62.4 |
| + remove input scaling (fix $\gamma_t \equiv 1$) | 63.6 |
| + operate on image input with inv. PCA | 63.6 |
| + operate on image output with inv. PCA | 63.9 |
| + predict original image | 64.5 |

Table 3. **Moving toward a classical DAE**, starting from our patch-wise PCA tokenizer. Each line is based on a modification of the immediately preceding line. See Sec. 4.3 for descriptions.

Interestingly, this *pixel-based* tokenizer exhibits a similar trend with other tokenizers we have studied, although the optimal dimension is shifted. In particular, the optimal dimension is $d$=48, which corresponds to an image size of 64 with a patch size of 4. With an image size of 256 and a patch size of 16 ($d$=768), the linear probe accuracy drops dramatically to 23.6%.

These comparisons show that the tokenizer and the resulting latent space are crucial for DDM/DAE to work competitively in the self-supervised learning scenario. In particular, applying a classical DAE with additive Gaussian noise on the pixel space leads to poor results.

### 4.3. Toward Classical Denoising Autoencoders

Next, we go on with our deconstruction trajectory and aim to get as close as possible to the classical DAE [39]. We attempt to remove every single aspect that still remains between our current PCA-based DDM and the classical DAE practice. Via this deconstructive process, we gain better understandings on how every modern design may influence the classical DAE. Tab. 3 gives the results, discussed next.

**Predict clean data (rather than noise).** While modern DDMs commonly predict the noise $\epsilon$ (see Eq. (2)), the classical DAE predicts the clean data instead. We examine this difference by minimizing the following loss function:

$$\lambda_t \| z_0 - \texttt{net}(z_t) \|^2 \qquad (3)$$

Here $z_0$ is the clean data (in the latent space), and $\texttt{net}(z_t)$ is the network prediction. $\lambda_t$ is a $t$-dependent loss weight, introduced to balance the contribution of different noise levels [34]. It is suggested to set $\lambda_t = \gamma_t^2/\sigma_t^2$ as per [34]. We find that setting $\lambda_t = \gamma_t^2$ works better in our scenario. Intuitively, it simply puts more weight to the loss terms of the *cleaner* data (larger $\gamma_t$).

With the modification of predicting clean data (rather than noise), the linear probe accuracy *degrades* from 65.1% to 62.4% (Tab. 3). This suggests that the choice of the prediction target influences the representation quality.

Even though we suffer from a degradation in this step, we will stick to this modification from now on, as our goal is to move towards a classical DAE.[4]

---

[4]We have revisited undoing this change in our final entry, in which we have not observed this degradation.

**Remove input scaling.** In modern DDMs (see Eq. (1)), the input is scaled by a factor of $\gamma_t$. This is not common practice in a classical DAE. Next, we study removing input scaling, *i.e.*, we set $\gamma_t \equiv 1$. As $\gamma_t$ is fixed, we need to define a noise schedule directly on $\sigma_t$. We simply set $\sigma_t$ as a linear schedule from 0 to $\sqrt{2}$. Moreover, we empirically set the weight in Eq. (3) as $\lambda_t = 1/(1 + \sigma_t^2)$, which again puts more emphasis on cleaner data (smaller $\sigma_t$).

After fixing $\gamma_t \equiv 1$, we achieve a decent accuracy of 63.6% (Tab. 3), which compares favorably with the varying $\gamma_t$ counterpart's 62.4%. This suggests that scaling the data by $\gamma_t$ is not necessary in our scenario.

**Operate on the *image* space with inverse PCA.** Thus far, for all entries we have explored (except Fig. 5), the model operates on the latent space produced by a tokenizer (Fig. 2 (b)). Ideally, we hope our DAE can work directly on the *image* space while still having good accuracy. With the usage of PCA, we can achieve this goal by inverse PCA.

The idea is illustrated in Fig. 1. Specially, we project the input image into the latent space by the PCA bases (*i.e.*, $V$), add noise in the latent space, and project the noisy latent back to the image space by the *inverse* PCA bases ($V^T$). Fig. 1 (middle, bottom) shows an example image with noise added in the latent space. With this noisy image as the input to the network, we can apply a standard ViT network [15] that directly operate on images, as if there is no tokenizer.

Applying this modification on the input side (while still predicting the output on the latent space) has 63.6% accuracy (Tab. 3). Further applying it to the output side (*i.e.*, predicting the output on the image space with inverse PCA) has 63.9% accuracy. Both results show that operating on the image space with inverse PCA can achieve similar results as operating on the latent space.

**Predict original image.** While inverse PCA can produce a prediction target in the image space, this target is not the original image. This is because PCA is a *lossy* encoder for any reduced dimension $d$. In contrast, it is a more natural solution to predict the original image directly.

When we let the network predict the original image, the "noise" introduced includes two parts: (i) the additive Gaussian noise, whose intrinsic dimension is $d$, and (ii) the PCA reconstruction error, whose intrinsic dimension is $D - d$ ($D$ is 768). We weight the loss of both parts differently.

Formally, with the clean original image $x_0$ and network prediction $\texttt{net}(x_t)$, we can compute the residue $r$ projected onto the full PCA space: $r \triangleq V(x_0 - \texttt{net}(x_t))$. Here $V$ is the $D$-by-$D$ matrix representing the full PCA bases. Then we minimize the following loss function:

$$\lambda_t \sum_{i=1}^{D} w_i r_i^2. \qquad (4)$$

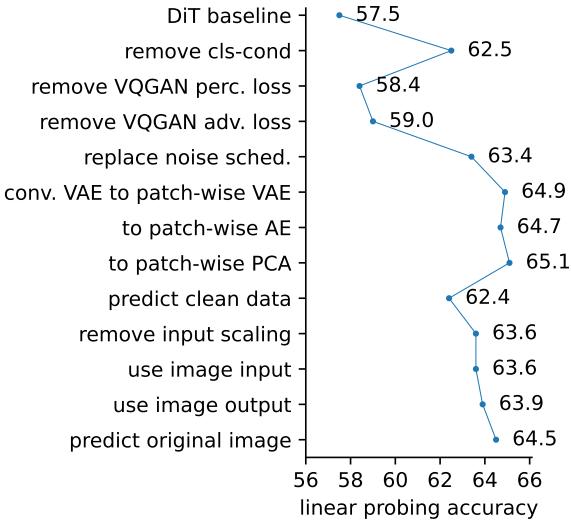Here $i$ denotes the $i$-th dimension of the vector $r$. The per-

Figure 6. **The overall deconstructive trajectory** from a modern DDM to *l*-DAE, summarizing Tab. 1, Tab. 2, and Tab. 3. Each line is based on a modification of the immediately preceding line.

dimension weight $w_i$ is 1 for $i \leq d$, and 0.1 for $d < i \leq D$. Intuitively, $w_i$ down-weights the loss of the PCA reconstruction error. With this formulation, predicting the original image achieves 64.5% linear probe accuracy (Tab. 3).

This variant is conceptually very simple: its input is a noisy image whose noise is added in the PCA latent space, its prediction is the original clean image (Fig. 1).

**Single noise level.** Lastly, out of curiosity, we further study a variant with *single-level* noise. We note that multi-level noise, given by noise scheduling, is a property motived by the diffusion process in DDMs; it is conceptually unnecessary in a classical DAE.

We fix the noise level $\sigma$ as a constant ($\sqrt{1/3}$). Using this single-level noise achieves decent accuracy of 61.5%, a 3% degradation *vs.* the multi-level noise counterpart (64.5%). Using multiple levels of noise is analogous to a form of data augmentation in DAE: it is beneficial, but not an enabling factor. This also implies that the representation capability of DDM is mainly gained by the denoising-driven process, not a diffusion-driven process.

As multi-level noise is useful and conceptually simple, we keep it in our final entries presented in the next section.

### 4.4. Summary

In sum, we deconstruct a modern DDM and push it towards a classical DAE (Fig. 6). We *undo* many of the modern designs and conceptually retain only two designs inherited from modern DDMs: (i) a low-dimensional latent space in which noise is added; and (ii) multi-level noise.

We use the entry at the end of Tab. 3 as our final DAE instantiation (illustrated in Fig. 1). We refer to this method as "*latent* Denoising Autoencoder", or in short, *l*-DAE.



Figure 7. **Visualization: pixel noise *vs*. latent noise. Left**: clean image, 256×256 pixels. **Middle**: Gaussian noise added to the pixel space. **Right**: Gaussian noise added to the latent space produced by the PCA tokenizer, visualized by back projection to the image space using inverse PCA. $\sigma = \sqrt{1/3}$ in both cases.

## 5. Analysis and Comparisons

**Visualizing latent noise.** Conceptually, *l*-DAE is a form of DAE that learns to remove noise added to the latent space. Thanks to the simplicity of PCA, we can easily visualize the latent noise by inverse PCA.

Fig. 7 compares the noise added to pixels *vs*. to the latent. Unlike the pixel noise, the latent noise is largely independent of the *resolution* of the image. With patch-wise PCA as the tokenizer, the pattern of the latent noise is mainly determined by the patch size. Intuitively, we may think of it as using patches, rather than pixels, to resolve the image. This behavior resembles MAE [21], which masks out patches instead of individual pixels.

**Denoising results.** Fig. 8 shows more examples of denoising results based on *l*-DAE. Our method produces reasonable predictions despite of the heavy noise. We note that this is less of a surprise, because neural network-based image restoration [4, 14] has been an intensively studied field.

Nevertheless, the visualization may help us better understand how *l*-DAE may learn good representations. The heavy noise added to the latent space creates a challenging *pretext* task for the model to solve. It is nontrivial (even for human beings) to predict the content based on one or a few noisy patches locally; the model is forced to learn higher-level, more holistic semantics to make sense of the underlying objects and scenes.

**Data augmentation.** Notably, all models we present thus far have *no data augmentation*: only the center crops of images are used, with no random resizing or color jittering, following [11, 32]. We further explore a mild data augmentation (random resized crop) for our final *l*-DAE:

| aug. | center crop | random crop |
|------|-------------|-------------|
| acc. | 64.5 | 65.0 |

which has slight improvement. This suggests that *the representation learning ability of l-DAE is largely independent of its reliance on data augmentation.* A similar behavior was observed in MAE [21], which sharply differs from the behavior of contrastive learning methods (*e.g.*, [6]).
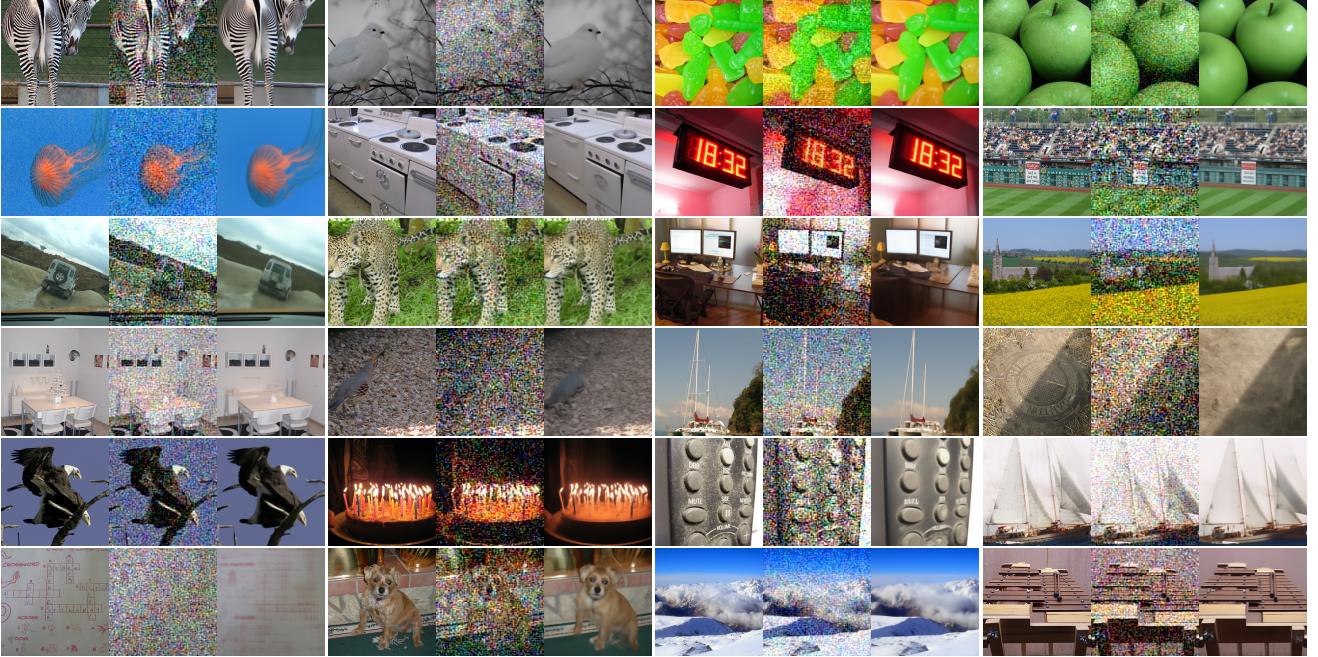
Figure 8. **Denoising results** of *l*-DAE, evaluated on ImageNet validation images. *This denoising problem, serving as a pretext task, encourages the network to learn meaningful representations in a self-supervised manner.* For each case, we show: (**left**) clean image; (**middle**) noisy image that is the input to the network, where the noise is added to the latent space; (**right**) denoised output.

**Training epochs.** All our experiments thus far are based on 400-epoch training. Following MAE [21], we also study training for 800 and 1600 epochs:

| epochs | 400 | 800 | 1600 |
|---|---|---|---|
| acc. | 65.0 | 67.5 | 69.6 |

As a reference, MAE [21] has a significant gain (4%) extending from 400 to 800 epochs, and MoCo v3 [7] has nearly no gain (0.2%) extending from 300 to 600 epochs.

**Model size.** Thus far, our all models are based on the DiT-L variant [32], whose encoder and decoder are both "ViT-$\frac{1}{2}$L" (half depth of ViT-L). We further train models of different sizes, whose encoder is ViT-B or ViT-L (decoder is always of the same size as encoder):

| encoder | ViT-B | ViT-$\frac{1}{2}$L | ViT-L |
|---|---|---|---|
| acc. | 60.3 | 65.0 | 70.9 |

We observe a good *scaling* behavior w.r.t. model size: scaling from ViT-B to ViT-L has a large gain of 10.6%. A similar scaling behavior is also observed in MAE [21], which has a 7.8% gain from ViT-B to ViT-L.

**Comparison with previous baselines.** Finally, to have a better sense of how different families of self-supervised learning methods perform, we compare with previous baselines in Tab. 4. We consider MoCo v3 [7], which belongs to the family of contrastive learning methods, and MAE [21], which belongs to the family of masking-based methods.

Interestingly, *l*-DAE performs decently in comparison with MAE, showing a degradation of 1.4% (ViT-B) or 0.8% (ViT-L). We note that here the training settings are made *as*

| method | ViT-B (86M) | ViT-L (304M) |
|---|---|---|
| MoCo v3 | 76.7 | 77.6 |
| MAE | 68.0 | 75.8 |
| *l*-DAE | 66.6 | 75.0 |

Table 4. **Comparisons with previous baselines** of MoCo v3 [7] and MAE [21]. The entries of both MAE and *l*-DAE are trained for 1600 epochs and with random crop as the data augmentation. Linear probe accuracy on ImageNet is reported. In the brackets are the number of parameters of the encoder.

*fair as possible* between MAE and *l*-DAE: both are trained for 1600 epochs and with random crop as the data augmentation. On the other hand, we should also note that MAE is more efficient in training because it only operates on unmasked patches. Nevertheless, *we have largely closed the accuracy gap between MAE and a DAE-driven method.*

Last, we observe that *autoencoder-based* methods (MAE and *l*-DAE) still fall short in comparison with contrastive learning methods under this protocol, especially when the model is small. We hope our study will draw more attention to the research on autoencoder-based methods for self-supervised learning.

## 6. Conclusion

We have reported that *l*-DAE, which largely resembles the classical DAE, can perform competitively in self-supervised learning. The critical component is a low-dimensional latent space on which noise is added. We hope our discovery will reignite interest in denoising-based methods in the context of today's self-supervised learning research.

8

## A. Implementation Details

**DiT architecture.** We follow the DiT architecture design [32]. The DiT architecture is similar to the original ViT [15], with extra modifications made for conditioning. Each Transformer block accepts an embedding network (a two-layer MLP) conditioned on the time step $t$. The output of this embedding network determines the scale and bias parameters of LayerNorm [1], referred to as adaLN [32]. Slightly different from [32], we set the hidden dimension of this MLP as 1/4 of its original dimension, which helps reduce model sizes and save memory, at no accuracy cost.

**Training.** The original DiTs [32] are trained with a batch size of 256. To speed up our exploration, we increase the batch size to 2048. We perform linear learning rate warm up [19] for 100 epochs and then decay it following a half-cycle cosine schedule. We use a base learning rate $blr$ = 1e-4 [32] by default, and set the actual $lr$ following the linear scaling rule [19]: $blr \times$ batch_size / 256. No weight decay is used [32]. We train for 400 epochs by default. On a 256-core TPU-v3 pod, training DiT-L takes 12 hours.

**Linear probing.** Our linear probing implementation follows the practice of MAE [21]. We use *clean*, 256×256-sized images for linear probing training and evaluation. The ViT output feature map is globally pooled by average pooling. It is then processed by a parameter-free BatchNorm [25] layer and a linear classifier layer, following [21]. The training batch size is 16384, learning rate is $6.4 \times 10^{-3}$ (cosine decay schedule), weight decay is 0, and training length is 90 epochs. Randomly resized crop and flipping are used during training and a single center crop is used for testing. Top-1 accuracy is reported.

While the model is conditioned on $t$ in self-supervised pre-training, conditioning is not needed in transfer learning (*e.g.*, linear probing). We fix the time step $t$ value in our linear probing training and evaluation. The influence of different $t$ values (out of 1000 time steps) is shown as follows:

| fixed $t$ | 0 | 10 | 20 | 40 | 80 |
|---|---|---|---|---|---|
| w/ clean input | 64.1 | 64.5 | 64.1 | 63.3 | 62.2 |
| w/ noisy input | 64.2 | 65.0 | 65.0 | 65.0 | 64.5 |

We note that the $t$ value determines: (i) the model weights, which are conditioned on $t$, and (ii) the noise added in transfer learning, using the same level of $t$. Both are shown in this table. We use $t = 10$ and clean input in all our experiments, except Tab. 4 where we use the optimal setting.

Fixing $t$ also means that the $t$-dependent MLP layers, which are used for conditioning, are not exposed in transfer learning, because they can be merged given the fixed $t$. As such, our model has the number of parameters just similar to the standard ViT [15], as reported in Tab. 4.

The DiT-L [32] has 24 blocks where the first 12 blocks are referred to as the "encoder" (hence ViT-$\frac{1}{2}$L) and the others the "decoder". This separation of the encoder and decoder is artificial. In the following table, we show the linear probing results using different numbers of blocks in the encoder, using the same pre-trained model:

| enc. blocks | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|
| acc. | 58.5 | 62.0 | 64.1 | **64.5** | 63.6 | 61.9 | 59.7 |

The optimal accuracy is achieved when the encoder and decoder have the same depth. This behavior is different from MAE's [21], whose encoder and decoder are asymmetric.

## B. Fine-tuning Results

In addition to linear probing, we also report end-to-end fine-tuning results. We closely followed MAE's protocol [21]. We use clean, 256×256-sized images as the inputs to the encoder. Globally average pooled outputs are used as features for classification. The training batch size is 1024, initial learning rate is $4 \times 10^{-3}$, weight decay is 0.05, drop path [24] is 0.1, and training length is 100 epochs. We use a layer-wise learning rate decay of 0.85 (B) or 0.65 (L). MixUp [43] (0.8), CutMix [42] (1.0), RandAug [8] (9, 0.5), and exponential moving average (0.9999) are used, similar to [21]. The results are summarized as below:

| method | ViT-B | ViT-L |
|---|---|---|
| MoCo v3 | 83.2 | 84.1 |
| MAE | 83.6 | 85.9 |
| *l*-DAE | 83.7 | 84.7 |

Overall, both autoencoder-based methods are better than MoCo v3. *l*-DAE performs similarly with MAE with ViT-B, but still fall short of MAE with ViT-L.

## References

[1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv:1607.06450*, 2016.

[2] Marcelo Bertalmio, Guillermo Sapiro, Vincent Caselles, and Coloma Ballester. Image inpainting. In *SIGGRAPH*, 2000.

[3] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*. Springer, 2006.

[4] Harold Christopher Burger, Christian J Schuler, and Stefan Harmeling. Image denoising: Can plain neural networks compete with BM3D? In *CVPR*, 2012.

[5] Kai Chen, Zhili Liu, Lanqing Hong, Hang Xu, Zhenguo Li, and Dit-Yan Yeung. Mixed autoencoder for self-supervised visual representation learning. In *CVPR*, 2023.

[6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020.

[7] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised Vision Transformers. In *ICCV*, 2021.

[8] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *CVPR Workshops*, 2020.

[9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009.

[10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2019.

[11] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat GANs on image synthesis. In *NeurIPS*, 2021.

[12] Jeff Donahue and Karen Simonyan. Large scale adversarial representation learning. *arXiv:1907.02544*, 2019.

[13] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. In *ICLR*, 2017.

[14] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image super-resolution. In *ECCV*, 2014.

[15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.

[16] Patrick Esser, Robin Rombach, and Björn Ommer. Taming transformers for high-resolution image synthesis. In *CVPR*, 2021.

[17] Yuxin Fang, Li Dong, Hangbo Bao, Xinggang Wang, and Furu Wei. Corrupted image modeling for self-supervised visual pre-training. *arXiv:2202.03382*, 2022.

[18] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, 2014.

[19] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch SGD: Training ImageNet in 1 hour. *arXiv:1706.02677*, 2017.

[20] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020.

[21] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, 2022.

[22] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local Nash equilibrium. In *NeurIPS*, 2017.

[23] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020.

[24] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *ECCV*, 2016.

[25] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.

[26] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *ICLR*, 2013.

[27] Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1989.

[28] Soumik Mukhopadhyay, Matthew Gwilliam, Vatsal Agarwal, Namitha Padmanabhan, Archana Swaminathan, Srinidhi Hegde, Tianyi Zhou, and Abhinav Shrivastava. Diffusion models beat GANs on image classification. *arXiv:2307.08702*, 2023.

[29] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *ICML*, 2021.

[30] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. In *NeurIPS*, 2017.

[31] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *CVPR*, 2016.

[32] William Peebles and Saining Xie. Scalable diffusion models with Transformers. In *ICCV*, 2023.

[33] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.

[34] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *ICLR*, 2022.

[35] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.

[36] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, 2015.

[37] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *NeurIPS*, 2019.

[38] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv:2011.13456*, 2020.

[39] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *ICML*, 2008.

[40] Weilai Xiang, Hongyu Yang, Di Huang, and Yunhong Wang. Denoising diffusion autoencoders are unified self-supervised learners. In *ICCV*, 2023.

[41] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. SimMIM: A simple framework for masked image modeling. In *CVPR*, 2022.

[42] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *ICCV*, 2019.

[43] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *ICLR*, 2018.

[44] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.