

High-Resolution Image Synthesis with Latent Diffusion Models

Robin Rombach¹ * Andreas Blattmann¹ * Dominik Lorenz¹ Patrick Esser^{RS} Björn Ommer¹

¹Ludwig Maximilian University of Munich & IWR, Heidelberg University, Germany ^{RS}Runway ML

<https://github.com/CompVis/latent-diffusion>

Abstract

By decomposing the image formation process into a sequential application of denoising autoencoders, diffusion models (DMs) achieve state-of-the-art synthesis results on image data and beyond. Additionally, their formulation allows for a guiding mechanism to control the image generation process without retraining. However, *since these models typically operate directly in pixel space, optimization of powerful DMs often consumes hundreds of GPU days and inference is expensive due to sequential evaluations.* *To enable DM training on limited computational resources while retaining their quality and flexibility,* we apply them in the latent space of powerful pretrained autoencoders. In contrast to previous work, training diffusion models on such a representation allows for the first time to reach a near-optimal point between complexity reduction and detail preservation, greatly boosting visual fidelity. By introducing cross-attention layers into the model architecture, we turn diffusion models into powerful and flexible generators for general conditioning inputs such as text or bounding boxes and high-resolution synthesis becomes possible in a convolutional manner. Our *latent diffusion models (LDMs)* achieve new state of the art scores for image inpainting and class-conditional image synthesis and highly competitive performance on various tasks, including unconditional image generation, text-to-image synthesis, and super-resolution, while significantly reducing computational requirements compared to pixel-based DMs.

1. Introduction

Image synthesis is one of the computer vision fields with the most spectacular recent development, but also among those with the greatest computational demands. Especially high-resolution synthesis of complex, natural scenes is presently dominated by scaling up likelihood-based models, potentially containing billions of parameters in autoregressive (AR) transformers [64, 65]. In contrast, the promising results of GANs [3, 26, 39] have been revealed to be mostly confined to data with comparably limited variability as their adversarial learning procedure does not easily scale to modeling complex, multi-modal distributions. Recently, diffusion models [79], which are built from a hierarchy of denoising autoencoders, have shown to achieve impressive



Figure 1. Boosting the upper bound on achievable quality with less aggressive downsampling. Since diffusion models offer excellent inductive biases for spatial data, we do not need the heavy spatial downsampling of related generative models in latent space, but can still greatly reduce the dimensionality of the data via suitable autoencoding models, see Sec. 3. Images are from the DIV2K [1] validation set, evaluated at 512^2 px. We denote the spatial downsampling factor by f . Reconstruction FIDs [28] and PSNR are calculated on ImageNet-val. [12]; see also Tab. 8.

results in image synthesis [29, 82] and beyond [7, 44, 47, 56], and define the state-of-the-art in class-conditional image synthesis [15, 30] and super-resolution [70]. Moreover, even unconditional DMs can readily be applied to tasks such as inpainting and colorization [82] or stroke-based synthesis [52], in contrast to other types of generative models [19, 45, 67]. Being likelihood-based models, they do not exhibit mode-collapse and training instabilities as GANs and, by heavily exploiting parameter sharing, they can model highly complex distributions of natural images without involving billions of parameters as in AR models [65].

Democratizing High-Resolution Image Synthesis DMs belong to the class of likelihood-based models, whose mode-covering behavior makes them prone to spend excessive amounts of capacity (and thus compute resources) on modeling imperceptible details of the data [16, 71]. Although the reweighted variational objective [29] aims to address this by undersampling the initial denoising steps, DMs are still computationally demanding, since training and evaluating such a model requires repeated function evaluations (and gradient computations) in the high-dimensional space of RGB images. As an example, training the most powerful DMs often takes hundreds of GPU days (e.g. 150 - 1000 V100 days in [15]) and repeated evaluations on a noisy version of the input space render also inference expensive,

*The first two authors contributed equally to this work.

so that producing 50k samples takes approximately 5 days [15] on a single A100 GPU. This has two consequences for the research community and users in general: Firstly, training such a model requires massive computational resources only available to a small fraction of the field, and leaves a huge carbon footprint [63, 83]. Secondly, evaluating an already trained model is also expensive in time and memory, since the same model architecture must run sequentially for a large number of steps (e.g. 25 - 1000 steps in [15]).

To increase the accessibility of this powerful model class and at the same time reduce its significant resource consumption, a method is needed that reduces the computational complexity for both training and sampling. Reducing the computational demands of DMs without impairing their performance is, therefore, key to enhance their accessibility.

Departure to Latent Space Our approach starts with the analysis of already trained diffusion models in pixel space: Fig. 2 shows the rate-distortion trade-off of a trained model. As with any likelihood-based model, learning can be roughly divided into two stages: First is a *perceptual compression* stage which removes high-frequency details but still learns little semantic variation. In the second stage, the actual generative model learns the semantic and conceptual composition of the data (*semantic compression*). We thus aim to first find a *perceptually equivalent, but computationally more suitable space*, in which we will train diffusion models for high-resolution image synthesis.

Following common practice [11, 23, 64, 65, 93], we separate training into two distinct phases: First, we train an autoencoder which provides a lower-dimensional (and thereby efficient) representational space which is perceptually equivalent to the data space. Importantly, and in contrast to previous work [23, 64], we do not need to rely on excessive spatial compression, as we train DMs in the learned latent space, which exhibits better scaling properties with respect to the spatial dimensionality. The reduced complexity also provides efficient image generation from the latent space with a single network pass. We dub the resulting model class *Latent Diffusion Models* (LDMs).

A notable advantage of this approach is that we need to train the universal autoencoding stage only once and can therefore reuse it for multiple DM trainings or to explore possibly completely different tasks [78]. This enables efficient exploration of a large number of diffusion models for various image-to-image and text-to-image tasks. For the latter, we design an architecture that connects transformers to the DM’s UNet backbone [69] and enables arbitrary types of token-based conditioning mechanisms, see Sec. 3.3.

In sum, our work makes the following **contributions**:

(i) In contrast to purely transformer-based approaches [23, 64], our method scales more gracefully to higher dimensional data and can thus (a) **work on a compression level which provides more faithful and detailed reconstructions than previous work** (see Fig. 1) and (b) can be efficiently

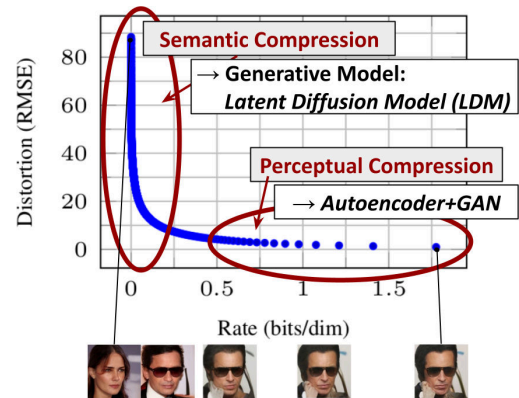


Figure 2. Illustrating perceptual and semantic compression: Most bits of a digital image correspond to imperceptible details. While DMs allow to suppress this semantically meaningless information by minimizing the responsible loss term, gradients (during training) and the neural network backbone (training and inference) still need to be evaluated on all pixels, leading to superfluous computations and unnecessarily expensive optimization and inference. We propose *latent diffusion models* (LDMs) as an effective generative model and a separate mild compression stage that only eliminates imperceptible details. Data and images from [29].

applied to high-resolution synthesis of megapixel images.

(ii) We achieve competitive performance on multiple tasks (unconditional image synthesis, inpainting, stochastic super-resolution) and datasets while significantly **lowering computational costs**. Compared to pixel-based diffusion approaches, we also significantly **decrease inference costs**.

(iii) We show that, in contrast to previous work [90] which learns both an encoder/decoder architecture and a score-based prior simultaneously, **our approach does not require a delicate weighting of reconstruction and generative abilities. This ensures extremely faithful reconstructions and requires very little regularization of the latent space.**

(iv) We find that **for densely conditioned tasks** such as super-resolution, inpainting and semantic synthesis, **our model can be applied in a convolutional fashion and render large, consistent images of $\sim 1024^2$ px.**

(v) Moreover, we **design a general-purpose conditioning mechanism based on cross-attention, enabling multi-modal training**. We use it to train class-conditional, text-to-image and layout-to-image models.

(vi) Finally, we release pretrained latent diffusion and autoencoding models at <https://github.com/CompVis/latent-diffusion> which might be reusable for a various tasks besides training of DMs [78].

2. Related Work

Generative Models for Image Synthesis The high dimensional nature of images presents distinct challenges to generative modeling. Generative Adversarial Networks (GAN) [26] allow for efficient sampling of high resolution images with good perceptual quality [3, 41], but are diffi-

cult to optimize [2, 27, 53] and struggle to capture the full data distribution [54]. In contrast, likelihood-based methods emphasize good density estimation which renders optimization more well-behaved. Variational autoencoders (VAE) [45] and flow-based models [18, 19] enable efficient synthesis of high resolution images [9, 43, 89], but sample quality is not on par with GANs. While autoregressive models (ARM) [6, 10, 91, 92] achieve strong performance in density estimation, computationally demanding architectures [94] and a sequential sampling process limit them to low resolution images. Because pixel based representations of images contain barely perceptible, high-frequency details [16, 71], maximum-likelihood training spends a disproportionate amount of capacity on modeling them, resulting in long training times. To scale to higher resolutions, several two-stage approaches [23, 65, 97, 99] use ARMs to model a compressed latent image space instead of raw pixels.

Recently, **Diffusion Probabilistic Models** (DM) [79], have achieved state-of-the-art results in density estimation [44] as well as in sample quality [15]. The generative power of these models stems from a natural fit to the inductive biases of image-like data when their underlying neural backbone is implemented as a UNet [15, 29, 69, 82]. The best synthesis quality is usually achieved when a reweighted objective [29] is used for training. In this case, the DM corresponds to a lossy compressor and allow to trade image quality for compression capabilities. Evaluating and optimizing these models in pixel space, however, **has the downside of low inference speed and very high training costs**. While the former can be partially addressed by advanced sampling strategies [46, 73, 81] and hierarchical approaches [30, 90], **training on high-resolution image data always requires to calculate expensive gradients**. We address both drawbacks with our proposed *LDMs*, which work on a compressed latent space of lower dimensionality. This renders training computationally cheaper and speeds up inference with almost no reduction in synthesis quality (see Fig. 1).

Two-Stage Image Synthesis To mitigate the shortcomings of individual generative approaches, a lot of research [11, 23, 65, 68, 97, 99] has gone into combining the strengths of different methods into more efficient and performant models via a two stage approach. VQ-VAEs [65, 97] use autoregressive models to learn an expressive prior over a discretized latent space. [64] extend this approach to text-to-image generation by learning a joint distribution over discretized image and text representations. More generally, [68] uses conditionally invertible networks to provide a generic transfer between latent spaces of diverse domains. Different from VQ-VAEs, VQGANs [23, 99] employ a first stage with an adversarial and perceptual objective to scale autoregressive transformers to larger images. However, the high compression rates required for feasible ARM training, which introduces billions of trainable parameters [23, 64], limit the overall performance of such ap-

proaches and less compression comes at the price of high computational cost [23, 64]. Our work prevents such trade-offs, as our proposed *LDMs* scale more gently to higher dimensional latent spaces due to their convolutional backbone. Thus, we are free to choose the level of compression which optimally mediates between learning a powerful first stage, without leaving too much perceptual compression up to the generative diffusion model while guaranteeing high-fidelity reconstructions (see Fig. 1). While approaches to jointly learn an encoding/decoding model together with a score-based prior exist [90], they still require a difficult weighting between reconstruction and generative capabilities [11] and are outperformed by our approach (Sec. 4).

3. Method

To lower the computational demands of training diffusion models towards high-resolution image synthesis, we observe that although diffusion models allow to ignore perceptually irrelevant details by undersampling the corresponding loss terms [29], they still require costly function evaluations in pixel space, which causes huge demands in computation time and energy resources.

We propose to circumvent this drawback by introducing an explicit separation of the compressive from the generative learning phase (see Fig. 2). To achieve this, we utilize an autoencoding model which learns a space that is perceptually equivalent to the image space, but offers significantly reduced computational complexity.

Such an approach offers several advantages: (i) By leaving the high-dimensional image space, we obtain DMs which are computationally much more efficient because sampling is performed on a low-dimensional space. (ii) We exploit the inductive bias of DMs inherited from their UNet architecture [69], which makes them particularly effective for data with spatial structure and therefore alleviates the need for aggressive, quality-reducing compression levels as required by previous approaches [23, 64]. (iii) Finally, we obtain general-purpose compression models whose latent space can be used to train multiple generative models and which can also be utilized for other downstream applications such as single-image CLIP-guided synthesis [25].

3.1. Perceptual Image Compression

Our perceptual compression model is based on previous work [23] and consists of an autoencoder trained by combination of **a perceptual loss** [102] and **a patch-based [32] adversarial objective** [20, 23, 99]. This ensures that the reconstructions are confined to the image manifold by enforcing local realism and avoids blurriness introduced by relying solely on pixel-space losses such as L_2 or L_1 objectives.

More precisely, given an image $x \in \mathbb{R}^{H \times W \times 3}$ in RGB space, the encoder \mathcal{E} encodes x into a latent representation $z = \mathcal{E}(x)$, and the decoder \mathcal{D} reconstructs the image from the latent, giving $\tilde{x} = \mathcal{D}(z) = \mathcal{D}(\mathcal{E}(x))$, where

$z \in \mathbb{R}^{h \times w \times c}$. Importantly, the encoder *downsamples* the image by a factor $f = H/h = W/w$, and we investigate different downsampling factors $f = 2^m$, with $m \in \mathbb{N}$.

In order to avoid arbitrarily high-variance latent spaces, we experiment with two different kinds of regularizations. The first variant, *KL-reg.*, imposes a slight KL-penalty towards a standard normal on the learned latent, similar to a VAE [45, 67], whereas *VQ-reg.* uses a vector quantization layer [93] within the decoder. This model can be interpreted as a VQGAN [23] but with the quantization layer absorbed by the decoder. Because our subsequent DM is designed to work with the two-dimensional structure of our learned latent space $z = \mathcal{E}(x)$, we can use relatively mild compression rates and achieve very good reconstructions. This is in contrast to previous works [23, 64], which relied on an arbitrary 1D ordering of the learned space z to model its distribution autoregressively and thereby ignored much of the inherent structure of z . Hence, our compression model preserves details of x better (see Tab. 8). The full objective and training details can be found in the supplement.

3.2. Latent Diffusion Models

Diffusion Models [79] are probabilistic models designed to learn a data distribution $p(x)$ by gradually denoising a normally distributed variable, which corresponds to learning the reverse process of a fixed Markov Chain of length T . For image synthesis, the most successful models [15, 29, 70] rely on a reweighted variant of the variational lower bound on $p(x)$, which mirrors denoising score-matching [82]. These models can be interpreted as an equally weighted sequence of denoising autoencoders $\epsilon_\theta(x_t, t)$; $t = 1 \dots T$, which are trained to predict a denoised variant of their input x_t , where x_t is a noisy version of the input x . The corresponding objective can be simplified to (Sec. A)

$$L_{DM} = \mathbb{E}_{x, \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_\theta(x_t, t)\|_2^2 \right], \quad (1)$$

with t uniformly sampled from $\{1, \dots, T\}$.

Generative Modeling of Latent Representations With our trained perceptual compression models consisting of \mathcal{E} and \mathcal{D} , we now have access to an efficient, low-dimensional latent space in which high-frequency, imperceptible details are abstracted away. Compared to the high-dimensional pixel space, this space is more suitable for likelihood-based generative models, as they can now (i) focus on the important, semantic bits of the data and (ii) train in a lower dimensional, computationally much more efficient space.

Unlike previous work that relied on autoregressive, attention-based transformer models in a highly compressed, discrete latent space [23, 64, 99], we can take advantage of image-specific inductive biases that our model offers. This includes the ability to build the underlying UNet primarily from 2D convolutional layers, and further focusing the

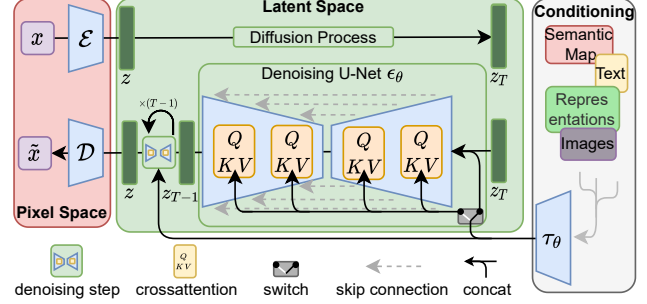


Figure 3. We condition LDMs either via concatenation or by a more general cross-attention mechanism. See Sec. 3.3

objective on the perceptually most relevant bits using the reweighted bound, which now reads

$$L_{LDM} := \mathbb{E}_{\mathcal{E}(x), \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_\theta(z_t, t)\|_2^2 \right]. \quad (2)$$

The neural backbone $\epsilon_\theta(\circ, t)$ of our model is realized as a time-conditional UNet [69]. Since the forward process is fixed, z_t can be efficiently obtained from \mathcal{E} during training, and samples from $p(z)$ can be decoded to image space with a single pass through \mathcal{D} .

3.3. Conditioning Mechanisms

Similar to other types of generative models [55, 80], diffusion models are in principle capable of modeling conditional distributions of the form $p(z|y)$. This can be implemented with a conditional denoising autoencoder $\epsilon_\theta(z_t, t, y)$ and paves the way to controlling the synthesis process through inputs y such as text [66], semantic maps [32, 59] or other image-to-image translation tasks [33].

In the context of image synthesis, however, combining the generative power of DMs with other types of conditionings beyond class-labels [15] or blurred variants of the input image [70] is so far an under-explored area of research.

We turn DMs into more flexible conditional image generators by augmenting their underlying UNet backbone with the cross-attention mechanism [94], which is effective for learning attention-based models of various input modalities [34, 35]. To pre-process y from various modalities (such as language prompts) we introduce a domain specific encoder τ_θ that projects y to an intermediate representation $\tau_\theta(y) \in \mathbb{R}^{M \times d_\tau}$, which is then mapped to the intermediate layers of the UNet via a cross-attention layer implementing $\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d}} \right) \cdot V$, with

$$Q = W_Q^{(i)} \cdot \varphi_i(z_t), \quad K = W_K^{(i)} \cdot \tau_\theta(y), \quad V = W_V^{(i)} \cdot \tau_\theta(y).$$

Here, $\varphi_i(z_t) \in \mathbb{R}^{N \times d_\epsilon^i}$ denotes a (flattened) intermediate representation of the UNet implementing ϵ_θ and $W_V^{(i)} \in \mathbb{R}^{d \times d_\epsilon^i}$, $W_Q^{(i)} \in \mathbb{R}^{d \times d_\tau}$ & $W_K^{(i)} \in \mathbb{R}^{d \times d_\tau}$ are learnable projection matrices [35, 94]. See Fig. 3 for a visual depiction.



Figure 4. Samples from *LDMs* trained on CelebAHQ [38], FFHQ [40], LSUN-Churches [98], LSUN-Bedrooms [98] and class-conditional ImageNet [12], each with a resolution of 256×256 . Best viewed when zoomed in. For more samples *cf.* the supplement.

Based on image-conditioning pairs, we then learn the conditional LDM via

$$L_{LDM} := \mathbb{E}_{\mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_{\theta}(z_t, t, \tau_{\theta}(y))\|_2^2 \right], \quad (3)$$

where both τ_{θ} and ϵ_{θ} are jointly optimized via Eq. 3. This conditioning mechanism is flexible as τ_{θ} can be parameterized with domain-specific experts, *e.g.* (unmasked) transformers [94] when y are text prompts (see Sec. 4.3.1)

4. Experiments

LDMs provide means to flexible and computationally tractable diffusion based image synthesis also including high-resolution generation of various image modalities, which we empirically show in the following. Firstly, however, we analyze the gains of our models compared to pixel-based diffusion models in both training and inference. Interestingly, we find that *LDMs* trained in *VQ*-regularized latent spaces achieve better sample quality, even though the reconstruction capabilities of *VQ*-regularized first stage models slightly fall behind those of their continuous counterparts, *cf.* Tab. 8. Therefore, we evaluate *VQ*-regularized *LDMs* in the remainder of the paper, unless stated differently. A visual comparison between the effects of first stage regularization schemes on *LDM* training and their generalization abilities to resolutions higher than 256^2 can be found in Appendix C.1. In D.2 we furthermore list details on architecture, implementation, training and evaluation for all results presented in this section.

4.1. On Perceptual Compression Tradeoffs

This section analyzes the behavior of our *LDMs* with different downsampling factors $f \in \{1, 2, 4, 8, 16, 32\}$ (abbreviated as *LDM-f*, where *LDM-1* corresponds to pixel-based DMs). To obtain a comparable test-field, we fix the computational resources to a single NVIDIA A100 for all experiments in this section and train all models for the same number of steps and with the same number of parameters.

Tab. 8 shows hyperparameters and reconstruction performance of the first stage models used for the *LDMs* compared in this section. Fig. 5 shows sample quality as a function of training progress for 2M steps of class-conditional

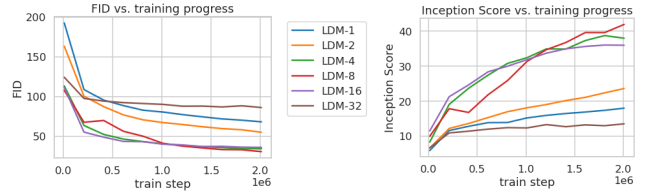


Figure 5. Analyzing the training of class-conditional *LDMs* with different downsampling factors f over 2M train steps on the ImageNet dataset. Pixel-based *LDM-1* requires substantially larger train times compared to models with larger downsampling factors (*LDM*-{4-16}). Too much perceptual compression as in *LDM-32* limits the overall sample quality. All models are trained on a single NVIDIA A100 with the same computational budget. Results obtained with 100 DDIM steps [81] and $\kappa = 0$.

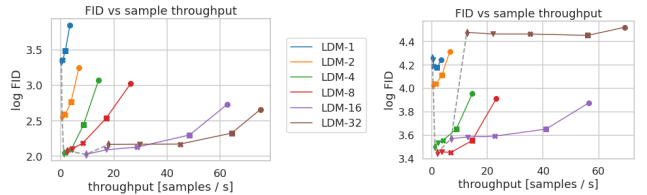


Figure 6. Inference speed vs sample quality: Comparing *LDMs* with different amounts of compression on the CelebA-HQ (left) and ImageNet (right) datasets. Different markers indicate $\{10, 20, 50, 100, 200\}$ sampling steps with the DDIM sampler, counted from right to left along each line. The dashed line shows the FID scores for 200 steps, indicating the strong performance of *LDM*-{4-8} compared to models with different compression ratios. FID scores assessed on 5000 samples. All models were trained for 500k (CelebA) / 2M (ImageNet) steps on an A100.

models on the ImageNet [12] dataset. We see that, i) small downsampling factors for *LDM*-{1,2} result in slow training progress, whereas ii) overly large values of f cause stagnating fidelity after comparably few training steps. Revisiting the analysis above (Fig. 1 and 2) we attribute this to i) leaving most of perceptual compression to the diffusion model and ii) too strong first stage compression resulting in information loss and thus limiting the achievable quality. *LDM*-{4-16} strike a good balance between efficiency and perceptually faithful results, which manifests in a sig-

CelebA-HQ 256 × 256					FFHQ 256 × 256				
Method	FID ↓	Prec. ↑	Recall ↑		Method	FID ↓	Prec. ↑	Recall ↑	
DC-VAE [61]	15.8	-	-		ImageBART [21]	9.57	-	-	
VQGAN+T [23] (k=400)	10.2	-	-		U-Net GAN (+aug) [75]	10.9 (7.6)	-	-	
PGGAN [38]	8.0	-	-		UDM [42]	5.54	-	-	
LSGM [90]	7.22	-	-		StyleGAN [40]	4.16	0.71	0.46	
UDM [42]	7.16	-	-		ProjectedGAN [74]	3.08	0.65	0.46	
<i>LDM-4</i> (ours, 500-s [†])	5.11	0.72	0.49		<i>LDM-4</i> (ours, 200-s)	4.98	0.73	0.50	

LSUN-Churches 256 × 256					LSUN-Bedrooms 256 × 256				
Method	FID ↓	Prec. ↑	Recall ↑		Method	FID ↓	Prec. ↑	Recall ↑	
DDPM [29]	7.89	-	-		ImageBART [21]	5.51	-	-	
ImageBART [21]	7.32	-	-		DDPM [29]	4.9	-	-	
PGGAN [38]	6.42	-	-		UDM [42]	4.57	-	-	
StyleGAN [40]	4.21	-	-		StyleGAN [40]	2.35	0.59	0.48	
StyleGAN2 [41]	3.86	-	-		ADM [15]	1.90	0.66	0.51	
ProjectedGAN [74]	1.59	0.61	0.44		ProjectedGAN [74]	1.52	0.61	0.34	
<i>LDM-8*</i> (ours, 200-s)	4.02	0.64	0.52		<i>LDM-4</i> (ours, 200-s)	2.95	0.66	0.48	

Table 1. Evaluation metrics for unconditional image synthesis. CelebA-HQ results reproduced from [42, 61, 96], FFHQ from [41, 42]. [†]: N -s refers to N sampling steps with the DDIM [81] sampler. ^{*}: trained in KL -regularized latent space. Additional results can be found in the supplementary.

nificant FID [28] gap of 38 between pixel-based diffusion ($LDM-1$) and $LDM-8$ after 2M training steps.

In Fig. 6, we compare models trained on CelebA-HQ [38] and ImageNet in terms sampling speed for different numbers of denoising steps with the DDIM sampler [81] and plot it against FID-scores [28]. $LDM-\{4-8\}$ outperform models with unsuitable ratios of perceptual and conceptual compression. Especially compared to pixel-based $LDM-1$, they achieve much lower FID scores while simultaneously significantly increasing sample throughput. Complex datasets such as ImageNet require reduced compression rates to avoid reducing quality. Summarized, we observe that $LDM-4$ and -8 lie in the best behaved regime for achieving high-quality synthesis results.

Text-Conditional Image Synthesis					
	DALL-E [†] [64]	CogView [†] [17]	Lafite [†] [105]	<i>LDM-KL-8</i>	<i>LDM-KL-8-G*</i>
FID ↓	27.50	27.10	26.94	23.35	12.61
IS ↑	17.90	18.20	26.02	19.93 \pm 0.35	26.62\pm0.38

Table 2. Evaluation of text-conditional image synthesis on the MS-COCO [50] dataset: Our model outperforms autoregressive [17, 64] and GAN-based [105] methods by a significant margin when using 250 DDIM [81] steps. [†]: Numbers taken from [105]. ^{*}: Classifier-free guidance [31], scale 1.5.

4.2. Image Generation with Latent Diffusion

We train unconditional models of 256^2 images on CelebA-HQ [38], FFHQ [40], LSUN-Churches and -Bedrooms [98] and evaluate the i) sample quality and ii) their coverage of the data manifold using ii) FID [28] and ii) Precision-and-Recall [49]. Tab. 1 summarizes our results. On CelebA-HQ, we report a new state-of-the-art FID of 5.11, outperforming previous likelihood-based models as well as GANs. We also outperform LSGM [90] where a latent diffusion model is trained jointly together with the first stage. In contrast, we train diffusion models in a fixed space and avoid the difficulty of weighing reconstruction quality against learning the prior over the latent space, see Fig. 1-2.

We outperform prior diffusion based approaches on all but the LSUN-Bedrooms dataset, where our score is close

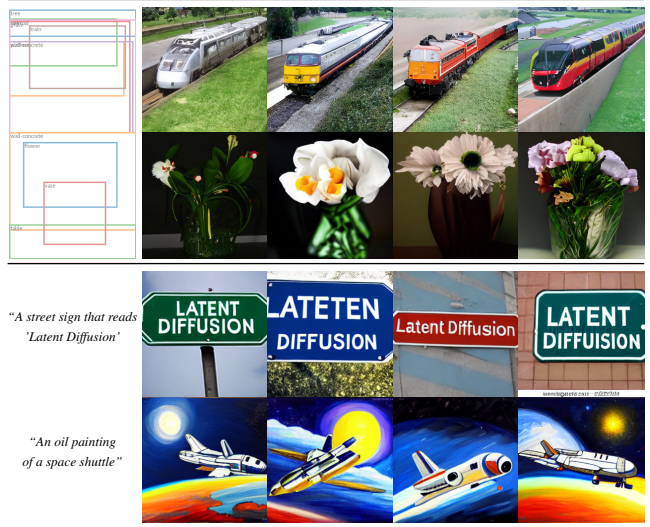


Figure 7. *Top*: Samples from our LDM for layout-to-image synthesis on COCO [4]. Quantitative evaluation in the supplement. *Bottom*: Samples from our text-to-image LDM model for user-defined text prompts, which is trained on LAION-400M [76].

to ADM [15], despite utilizing half its parameters and requiring 4-times less train resources (see Appendix D.3.5). Moreover, $LDMs$ consistently improve upon GAN-based methods in Precision and Recall, thus confirming the advantages of their mode-covering likelihood-based training objective over adversarial approaches. In Fig. 4 we also show qualitative results on each dataset.

4.3. Conditional Latent Diffusion

4.3.1 Transformer Encoders for LDMs

By introducing cross-attention based conditioning into LDMs we open them up for various conditioning modalities previously unexplored for diffusion models. For **text-to-image** image modeling, we train a 1.45B parameter model conditioned on language prompts on LAION-400M [76]. We employ the BERT-tokenizer [14] and implement τ_θ as a transformer [94] to infer a latent code which is mapped into the UNet via cross-attention (Sec. 3.3). This combination of domain specific experts for learning a language representation and visual synthesis results in a powerful model, which generalizes well to complex, user-defined text prompts, cf. Fig. 7 and 14. For quantitative analysis, we follow prior work and evaluate text-to-image generation on the MS-COCO [50] validation set, where our model improves upon powerful AR [17, 64] and GAN-based [105] methods, cf. Tab. 2. We note that applying classifier-free diffusion guidance [31] greatly boosts sample quality. To further analyze the flexibility of the cross-attention based conditioning mechanism we also train models to synthesize images based on **semantic layouts** on OpenImages [48], and finetune on COCO [4], see Fig. 7. See Sec. C.4 for the quantitative evaluation and implementation details.

Lastly, following prior work [3, 15, 21, 23], we evaluate our best-performing class-conditional ImageNet mod-

els with $f \in \{4, 8\}$ from Sec. 4.1 in Tab. 3, Fig. 4 and Sec. C.5. Here we outperform the state of the art diffusion model ADM [15] while significantly reducing computational requirements and parameter count, cf. Tab 18.

Method	FID↓	IS↑	Precision↑	Recall↑	N_{params}	
BigGan-deep [3]	6.95	203.6 \pm 2.6	0.87	0.28	340M	
ADM [15]	10.94	100.98	0.69	0.63	554M	250 DDIM steps
ADM-G [15]	4.59	186.7	0.82	0.52	608M	250 DDIM steps
LDM-4 (ours)	10.56	103.49 \pm 1.24	0.71	0.62	400M	250 DDIM steps
LDM-4-G (ours)	3.60	247.67\pm5.09	0.87	0.48	400M	250 steps, classifier-free guidance [31], scale 1.5

Table 3. Comparison of a class-conditional ImageNet LDM with recent state-of-the-art methods for class-conditional image generation on the ImageNet [12] dataset. A more detailed comparison with additional baselines can be found in C.5, Tab. 10 and E.

4.3.2 Convolutional Sampling Beyond 256²

By concatenating spatially aligned conditioning information to the input of ϵ_θ , LDMs can serve as efficient general-purpose image-to-image translation models. We use this to train models for semantic synthesis, super-resolution (Sec. 4.4) and inpainting (Sec. 4.5). For semantic synthesis, we use images of landscapes paired with semantic maps [23, 59] and concatenate downsampled versions of the semantic maps with the latent image representation of a $f = 4$ model (VQ-reg., see Tab. 8). We train on an input resolution of 256² (crops from 384²) but find that our model generalizes to larger resolutions and can generate images up to the megapixel regime when evaluated in a convolutional manner (see Fig. 8). We exploit this behavior to also apply the super-resolution models in Sec. 4.4 and the inpainting models in Sec. 4.5 to generate large images between 512² and 1024². For this application, the signal-to-noise ratio (induced by the scale of the latent space) significantly affects the results. In Sec. C.1 we illustrate this when learning an LDM on (i) the latent space as provided by a $f = 4$ model (KL-reg., see Tab. 8), and (ii) a rescaled version, scaled by the component-wise standard deviation.



Figure 8. A LDM trained on 256² resolution can generalize to larger resolution (here: 512 × 1024) for spatially conditioned tasks such as semantic synthesis of landscape images. See Sec. 4.3.2.

4.4. Super-Resolution with Latent Diffusion

LDMs can be efficiently trained for super-resolution by directly conditioning on low-resolution images via concate-



Figure 9. ImageNet 64→256 super-resolution on ImageNet-Val. LDM-SR has advantages at rendering realistic textures but SR3 can synthesize more coherent fine structures. See appendix for additional samples and cropouts. SR3 results from [70].

nation (cf. Sec. 3.3). In a first experiment, we follow SR3 [70] and fix the image degradation to a bicubic interpolation with 4×-downsampling and train on ImageNet following SR3’s data processing pipeline. We use the $f = 4$ autoencoding model pretrained on OpenImages (VQ-reg., cf. Tab. 8) and concatenate the low-resolution conditioning y and the inputs to the UNet, i.e. τ_θ is the identity. Our qualitative and quantitative results (see Fig. 9 and Tab. 4) show competitive performance and LDM-SR outperforms SR3 in FID while SR3 has a better IS. A simple image regression model achieves the highest PSNR and SSIM scores; however these metrics do not align well with human perception [102] and favor blurriness over imperfectly aligned high frequency details [70]. Further, we conduct a user study comparing the pixel-baseline with LDM-SR. We follow SR3 [70] where human subjects were shown a low-res image in between two high-res images and asked for preference. The results in Tab. 5 affirm the good performance of LDM-SR. PSNR and SSIM can be pushed by using a post-hoc guiding mechanism [15] and we implement this *image-based guider* via a perceptual loss, see Sec. C.7. Since the bicubic degradation process does not generalize well to images which do not follow this pre-processing, we also train a generic model, LDM-BSR, by using more diverse degradation. The results are shown in Sec. C.7.1.

Method	FID ↓	IS ↑	PSNR ↑	SSIM ↑	N_{params}	Throughput* (samples/s)
Image Regression [70]	15.2	121.1	27.9	0.801	625M	N/A
SR3 [70]	5.2	180.1	26.4	0.762	625M	N/A
LDM-4 (ours, 100 steps)	2.8 [†] /4.8 [‡]	166.3	24.4 \pm 3.8	0.69 \pm 0.14	169M	4.62
LDM-4 (ours, big, 100 steps)	2.4[†]/4.3[‡]	174.9	24.7 \pm 4.1	0.71 \pm 0.15	552M	4.5
LDM-4 (ours, 50 steps, guiding)	4.4 [†] /6.4 [‡]	153.7	25.8 \pm 3.7	0.74 \pm 0.12	184M	0.38

Table 4. ×4 upscaling results on ImageNet-Val. (256²); [†]: FID features computed on validation split, [‡]: FID features computed on train split; *: Assessed on a NVIDIA A100

4.5. Inpainting with Latent Diffusion

Inpainting is the task of filling masked regions of an image with new content either because parts of the image are

User Study	SR on ImageNet		Inpainting on Places	
	Pixel-DM (f_1)	<i>LDM-4</i>	LAMA [85]	<i>LDM-4</i>
Task 1: Preference vs GT \uparrow	16.0%	30.4%	13.6%	21.0%
Task 2: Preference Score \uparrow	29.4%	70.6%	31.9%	68.1%

Table 5. Task 1: Subjects were shown ground truth and generated image and asked for preference. Task 2: Subjects had to decide between two generated images. More details in D.3.6

Model (reg.-type)	train throughput samples/sec.	sampling throughput [†] @256	@512	train+val hours/epoch	FID@2k epoch 6
<i>LDM-1</i> (no first stage)	0.11	0.26	0.07	20.66	24.74
<i>LDM-4</i> (KL, w/ attn)	0.32	0.97	0.34	7.66	15.21
<i>LDM-4</i> (VQ, w/ attn)	0.33	0.97	0.34	7.04	14.99
<i>LDM-4</i> (VQ, w/o attn)	0.35	0.99	0.36	6.66	15.95

Table 6. Assessing inpainting efficiency. [†]: Deviations from Fig. 6 due to varying GPU settings/batch sizes *cf.* the supplement.

are corrupted or to replace existing but undesired content within the image. We evaluate how our general approach for conditional image generation compares to more specialized, state-of-the-art approaches for this task. Our evaluation follows the protocol of LaMa [85], a recent inpainting model that introduces a specialized architecture relying on Fast Fourier Convolutions [8]. We describe the exact training & evaluation protocol on Places [104] in Sec. D.2.2.

We first analyze the effect of different design choices for the first stage. We compare the inpainting efficiency of *LDM-1* (*i.e.* a pixel-based conditional DM) with *LDM-4*, for both KL and VQ regularizations, as well as *VQ-LDM-4* without any attention in the first stage (see Tab. 8), where the latter reduces GPU memory for decoding at high resolutions. For comparability, we fix the number of parameters for all models. Tab. 6 reports the training and sampling throughput at resolution 256^2 and 512^2 , the total training time in hours per epoch and the FID score on the validation split after six epochs. Overall, we observe a speed-up of at least $2.7\times$ between pixel- and latent-based diffusion models while improving FID scores by a factor of at least $1.6\times$.

The comparison with other inpainting approaches in Tab. 7 shows that our model with attention improves the overall image quality as measured by FID over that of [85]. LPIPS between the unmasked images and our samples is slightly higher than that of [85]. We attribute this to [85] only producing a single result which tends to recover more of an average image compared to the diverse results produced by our LDM *cf.* Fig. 20. Additionally in a user study (Tab. 5) human subjects favor our results over those of [85].

Based on these initial results, we also trained a larger diffusion model (*big* in Tab. 7) in the latent space of the VQ-regularized first stage without attention. Following [15], the UNet of this diffusion model uses attention layers on three levels of its feature hierarchy, the BigGAN [3] residual block for up- and downsampling and has 387M parameters instead of 215M. After training, we noticed a discrepancy in the quality of samples produced at resolutions 256^2 and 512^2 , which we hypothesize to be caused by the additional attention modules. However, fine-tuning the model for half

Method	40-50% masked		All samples	
	FID \downarrow	LPIPS \downarrow	FID \downarrow	LPIPS \downarrow
<i>LDM-4</i> (ours, big, w/ ft)	9.39	0.246 ± 0.042	1.50	0.137 ± 0.080
<i>LDM-4</i> (ours, big, w/o ft)	12.89	0.257 ± 0.047	2.40	0.142 ± 0.085
<i>LDM-4</i> (ours, w/ attn)	11.87	0.257 ± 0.042	2.15	0.144 ± 0.084
<i>LDM-4</i> (ours, w/o attn)	12.60	0.259 ± 0.041	2.37	0.145 ± 0.084
LaMa [85] [†]	12.31	0.243 ± 0.038	2.23	0.134 ± 0.080
LaMa [85]	12.0	0.24	2.21	0.14
CoModGAN [103]	10.4	0.26	<u>1.82</u>	0.15
RegionWise [51]	21.3	0.27	4.75	0.15
DeepFill v2 [100]	22.1	0.28	5.20	0.16
EdgeConnect [57]	30.5	0.28	8.37	0.16

Table 7. Comparison of inpainting performance on 30k crops of size 512×512 from test images of Places [104]. The column 40-50% reports metrics computed over hard examples where 40-50% of the image region have to be inpainted. [†]recomputed on our test set, since the original test set used in [85] was not available.

an epoch at resolution 512^2 allows the model to adjust to the new feature statistics and sets a new state of the art FID on image inpainting (*big, w/o attn, w/ ft* in Tab. 7, Fig. 10.).

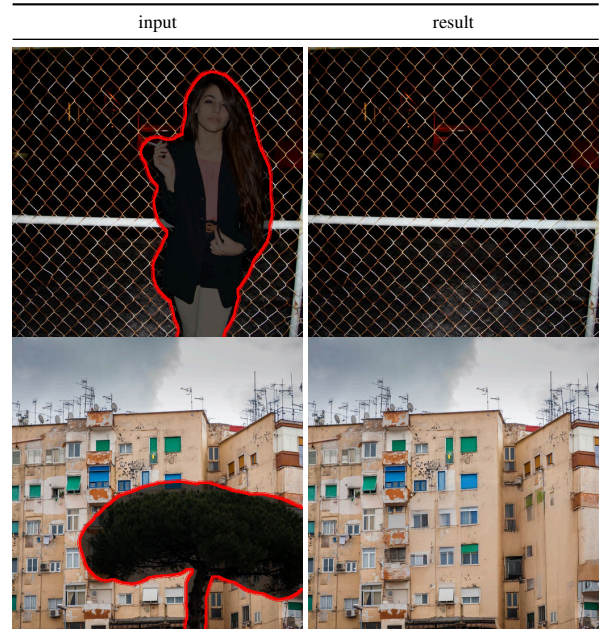


Figure 10. Qualitative results on object removal with our *big, w/ ft* inpainting model. For more results, see Fig. 21.

5. Conclusion

We have presented latent diffusion models, a simple and efficient way to significantly improve both the training and sampling efficiency of denoising diffusion models without degrading their quality. Based on this and our cross-attention conditioning mechanism, our experiments could demonstrate favorable results compared to state-of-the-art methods across a wide range of conditional image synthesis tasks without task-specific architectures.

This work has been supported by the German Federal Ministry for Economic Affairs and Energy within the project KI-Absicherung - Safe AI for automated driving and by the German Research Foundation (DFG) project 421703927.

References

- [1] Eirikur Agustsson and Radu Timofte. NTIRE 2017 challenge on single image super-resolution: Dataset and study. In *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 1122–1131. IEEE Computer Society, 2017. 1
- [2] Martin Arjovsky, Soumith Chintala, and Lon Bottou. Wasserstein gan, 2017. 3
- [3] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *Int. Conf. Learn. Represent.*, 2019. 1, 2, 6, 7, 8, 19, 26
- [4] Holger Caesar, Jasper R. R. Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 1209–1218. Computer Vision Foundation / IEEE Computer Society, 2018. 6, 17, 18
- [5] Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650, 2021. 27
- [6] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pre-training from pixels. In *ICML*, volume 119 of *Proceedings of Machine Learning Research*, pages 1691–1703. PMLR, 2020. 3
- [7] Nanxin Chen, Yu Zhang, Heiga Zen, Ron J. Weiss, Mohammad Norouzi, and William Chan. Wavegrad: Estimating gradients for waveform generation. In *ICLR*. OpenReview.net, 2021. 1
- [8] Lu Chi, Borui Jiang, and Yadong Mu. Fast fourier convolution. In *NeurIPS*, 2020. 8
- [9] Rewon Child. Very deep vaes generalize autoregressive models and can outperform them on images. *CoRR*, abs/2011.10650, 2020. 3
- [10] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *CoRR*, abs/1904.10509, 2019. 3
- [11] Bin Dai and David P. Wipf. Diagnosing and enhancing VAE models. In *ICLR (Poster)*. OpenReview.net, 2019. 2, 3
- [12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255. IEEE Computer Society, 2009. 1, 5, 7, 19
- [13] Emily Denton. Ethical considerations of generative ai. AI for Content Creation Workshop, CVPR, 2021. 27
- [14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018. 6
- [15] Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. *CoRR*, abs/2105.05233, 2021. 1, 2, 3, 4, 6, 7, 8, 15, 19, 23, 24, 26
- [16] Sander Dieleman. Musings on typicality, 2020. 1, 3
- [17] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, and Jie Tang. Cogview: Mastering text-to-image generation via transformers. *CoRR*, abs/2105.13290, 2021. 6
- [18] Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Non-linear independent components estimation, 2015. 3
- [19] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. 1, 3
- [20] Alexey Dosovitskiy and Thomas Brox. Generating images with perceptual similarity metrics based on deep networks. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett, editors, *Adv. Neural Inform. Process. Syst.*, pages 658–666, 2016. 3
- [21] Patrick Esser, Robin Rombach, Andreas Blattmann, and Björn Ommer. Imagebart: Bidirectional context with multinomial diffusion for autoregressive image synthesis. *CoRR*, abs/2108.08827, 2021. 6, 19
- [22] Patrick Esser, Robin Rombach, and Björn Ommer. A note on data biases in generative models. *arXiv preprint arXiv:2012.02516*, 2020. 27
- [23] Patrick Esser, Robin Rombach, and Björn Ommer. Taming transformers for high-resolution image synthesis. *CoRR*, abs/2012.09841, 2020. 2, 3, 4, 6, 7, 18, 19, 27, 32, 34
- [24] Mary Anne Franks and Ari Ezra Waldman. Sex, lies, and videotape: Deep fakes and free speech delusions. *Md. L. Rev.*, 78:892, 2018. 27
- [25] Kevin Frans, Lisa B. Soros, and Olaf Witkowski. Clipdraw: Exploring text-to-drawing synthesis through language-image encoders. *ArXiv*, abs/2106.14843, 2021. 3
- [26] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial networks. *CoRR*, 2014. 1, 2
- [27] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. Improved training of wasserstein gans, 2017. 3
- [28] 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 *Adv. Neural Inform. Process. Syst.*, pages 6626–6637, 2017. 1, 6, 24
- [29] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 1, 2, 3, 4, 6, 14
- [30] Jonathan Ho, Chitwan Saharia, William Chan, David J. Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *CoRR*, abs/2106.15282, 2021. 1, 3, 19
- [31] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021. 6, 7, 19, 26, 35, 36

- [32] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, pages 5967–5976. IEEE Computer Society, 2017. 3, 4
- [33] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation with conditional adversarial networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5967–5976, 2017. 4
- [34] Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David Ding, Skanda Koppula, Daniel Zoran, Andrew Brock, Evan Shelhamer, Olivier J. Hénaff, Matthew M. Botvinick, Andrew Zisserman, Oriol Vinyals, and João Carreira. Perceiver IO: A general architecture for structured inputs & outputs. *CoRR*, abs/2107.14795, 2021. 4
- [35] Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and João Carreira. Perceiver: General perception with iterative attention. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 4651–4664. PMLR, 2021. 4
- [36] Manuel Jahn, Robin Rombach, and Björn Ommer. High-resolution complex scene synthesis with transformers. *CoRR*, abs/2105.06458, 2021. 17, 18, 25
- [37] Niharika Jain, Alberto Olmo, Sailik Sengupta, Lydia Manikonda, and Subbarao Kambhampati. Imperfect imagination: Implications of gans exacerbating biases on facial data augmentation and snapchat selfie lenses. *arXiv preprint arXiv:2001.09528*, 2020. 27
- [38] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *CoRR*, abs/1710.10196, 2017. 5, 6
- [39] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 4401–4410, 2019. 1
- [40] T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial networks. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 5, 6
- [41] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. *CoRR*, abs/1912.04958, 2019. 2, 6, 26
- [42] Dongjun Kim, Seungjae Shin, Kyungwoo Song, Wanmo Kang, and Il-Chul Moon. Score matching model for unbounded data score. *CoRR*, abs/2106.05527, 2021. 6
- [43] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, 2018. 3
- [44] Diederik P. Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *CoRR*, abs/2107.00630, 2021. 1, 3, 14
- [45] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *2nd International Conference on Learning Representations, ICLR*, 2014. 1, 3, 4, 27
- [46] Zhifeng Kong and Wei Ping. On fast sampling of diffusion probabilistic models. *CoRR*, abs/2106.00132, 2021. 3
- [47] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. In *ICLR*. OpenReview.net, 2021. 1
- [48] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper R. R. Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Tom Duerig, and Vittorio Ferrari. The open images dataset V4: unified image classification, object detection, and visual relationship detection at scale. *CoRR*, abs/1811.00982, 2018. 6, 17, 18
- [49] Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved precision and recall metric for assessing generative models. *CoRR*, abs/1904.06991, 2019. 6, 24
- [50] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014. 6, 25
- [51] Yuqing Ma, Xianglong Liu, Shihao Bai, Le-Yi Wang, Aishan Liu, Dacheng Tao, and Edwin Hancock. Region-wise generative adversarial image inpainting for large missing areas. *ArXiv*, abs/1909.12507, 2019. 8
- [52] Chenlin Meng, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Image synthesis and editing with stochastic differential equations. *CoRR*, abs/2108.01073, 2021. 1
- [53] Lars M. Mescheder. On the convergence properties of GAN training. *CoRR*, abs/1801.04406, 2018. 3
- [54] Luke Metz, Ben Poole, David Pfau, and Jascha Sohl-Dickstein. Unrolled generative adversarial networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. 3
- [55] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *CoRR*, abs/1411.1784, 2014. 4
- [56] Gautam Mittal, Jesse H. Engel, Curtis Hawthorne, and Ian Simon. Symbolic music generation with diffusion models. *CoRR*, abs/2103.16091, 2021. 1
- [57] Kamyar Nazeri, Eric Ng, Tony Joseph, Faisal Z. Qureshi, and Mehran Ebrahimi. Edgeconnect: Generative image inpainting with adversarial edge learning. *ArXiv*, abs/1901.00212, 2019. 8
- [58] Anton Obukhov, Maximilian Seitzer, Po-Wei Wu, Semen Zhydenko, Jonathan Kyl, and Elvis Yu-Jing Lin. High-fidelity performance metrics for generative models in pytorch, 2020. Version: 0.3.0, DOI: 10.5281/zenodo.4957738. 24, 25
- [59] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 4, 7

- [60] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. **18**
- [61] Gaurav Parmar, Dacheng Li, Kwonjoon Lee, and Zhuowen Tu. Dual contradistinctive generative autoencoder. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pages 823–832. Computer Vision Foundation / IEEE, 2021. **6**
- [62] Gaurav Parmar, Richard Zhang, and Jun-Yan Zhu. On buggy resizing libraries and surprising subtleties in fid calculation. *arXiv preprint arXiv:2104.11222*, 2021. **24**
- [63] David A. Patterson, Joseph Gonzalez, Quoc V. Le, Chen Liang, Lluís-Miquel Munguia, Daniel Rothchild, David R. So, Maud Texier, and Jeff Dean. Carbon emissions and large neural network training. *CoRR*, abs/2104.10350, 2021. **2**
- [64] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. *CoRR*, abs/2102.12092, 2021. **1, 2, 3, 4, 6, 18, 25**
- [65] Ali Razavi, Aaron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with VQ-VAE-2. In *NeurIPS*, pages 14837–14847, 2019. **1, 2, 3, 19**
- [66] Scott E. Reed, Zeynep Akata, Xinchun Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. In *ICML*, 2016. **4**
- [67] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31st International Conference on International Conference on Machine Learning, ICML, 2014*. **1, 4, 27**
- [68] Robin Rombach, Patrick Esser, and Björn Ommer. Network-to-network translation with conditional invertible neural networks. In *NeurIPS*, 2020. **3**
- [69] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI (3)*, volume 9351 of *Lecture Notes in Computer Science*, pages 234–241. Springer, 2015. **2, 3, 4**
- [70] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J. Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *CoRR*, abs/2104.07636, 2021. **1, 4, 7, 19, 20, 21, 25**
- [71] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P. Kingma. Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications. *CoRR*, abs/1701.05517, 2017. **1, 3**
- [72] Dave Salvator. NVIDIA Developer Blog. <https://developer.nvidia.com/blog/getting-immediate-speedups-with-a100-tf32>, 2020. **26**
- [73] Robin San-Roman, Eliya Nachmani, and Lior Wolf. Noise estimation for generative diffusion models. *CoRR*, abs/2104.02600, 2021. **3**
- [74] Axel Sauer, Kashyap Chitta, Jens Müller, and Andreas Geiger. Projected gans converge faster. *CoRR*, abs/2111.01007, 2021. **6**
- [75] Edgar Schönfeld, Bernt Schiele, and Anna Khoreva. A unet based discriminator for generative adversarial networks. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 8204–8213. Computer Vision Foundation / IEEE, 2020. **6**
- [76] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs, 2021. **6**
- [77] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In Yoshua Bengio and Yann LeCun, editors, *Int. Conf. Learn. Represent.*, 2015. **27, 41, 42, 43**
- [78] Charlie Snell. Alien Dreams: An Emerging Art Scene. <https://ml.berkeley.edu/blog/posts/clip-art/>, 2021. [Online; accessed November-2021]. **2**
- [79] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. *CoRR*, abs/1503.03585, 2015. **1, 3, 4, 15**
- [80] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. **4**
- [81] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*. OpenReview.net, 2021. **3, 5, 6, 20**
- [82] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *CoRR*, abs/2011.13456, 2020. **1, 3, 4, 15**
- [83] Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and policy considerations for modern deep learning research. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 13693–13696. AAAI Press, 2020. **2**
- [84] Wei Sun and Tianfu Wu. Learning layout and style reconfigurable gans for controllable image synthesis. *CoRR*, abs/2003.11571, 2020. **18, 25**
- [85] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor S. Lempitsky. Resolution-robust large mask inpainting with fourier convolutions. *ArXiv*, abs/2109.07161, 2021. **8, 24, 30**
- [86] Tristan Sylvain, Pengchuan Zhang, Yoshua Bengio, R. De-von Hjelm, and Shikhar Sharma. Object-centric image generation from layouts. In *Thirty-Fifth AAAI Conference on*

- Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 2647–2655. AAAI Press, 2021. 17, 18, 25
- [87] Patrick Tinsley, Adam Czajka, and Patrick Flynn. This face does not exist... but it might be yours! identity leakage in generative models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1320–1328, 2021. 27
- [88] Antonio Torralba and Alexei A Efros. Unbiased look at dataset bias. In *CVPR 2011*, pages 1521–1528. IEEE, 2011. 27
- [89] Arash Vahdat and Jan Kautz. NVAE: A deep hierarchical variational autoencoder. In *NeurIPS*, 2020. 3
- [90] Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space. *CoRR*, abs/2106.05931, 2021. 2, 3, 6
- [91] Aaron van den Oord, Nal Kalchbrenner, Lasse Espeholt, koray kavukcuoglu, Oriol Vinyals, and Alex Graves. Conditional image generation with pixelcnn decoders. In *Advances in Neural Information Processing Systems*, 2016. 3
- [92] Aäron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. *CoRR*, abs/1601.06759, 2016. 3
- [93] Aäron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. In *NIPS*, pages 6306–6315, 2017. 2, 4, 27
- [94] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, pages 5998–6008, 2017. 3, 4, 5, 6
- [95] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. Huggingface’s transformers: State-of-the-art natural language processing. *CoRR*, abs/1910.03771, 2019. 24
- [96] Zhisheng Xiao, Karsten Kreis, Jan Kautz, and Arash Vahdat. VAEBM: A symbiosis between variational autoencoders and energy-based models. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. 6
- [97] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. Videogpt: Video generation using VQ-VAE and transformers. *CoRR*, abs/2104.10157, 2021. 3
- [98] Fisher Yu, Yinda Zhang, Shuran Song, Ari Seff, and Jianxiong Xiao. LSUN: construction of a large-scale image dataset using deep learning with humans in the loop. *CoRR*, abs/1506.03365, 2015. 5, 6
- [99] Jiahui Yu, Xin Li, Jing Yu Koh, Han Zhang, Ruoming Pang, James Qin, Alexander Ku, Yuanzhong Xu, Jason Baldridge, and Yonghui Wu. Vector-quantized image modeling with improved vqgan, 2021. 3, 4
- [100] Jiahui Yu, Zhe L. Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S. Huang. Free-form image inpainting with gated convolution. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4470–4479, 2019. 8
- [101] K. Zhang, Jingyun Liang, Luc Van Gool, and Radu Timofte. Designing a practical degradation model for deep blind image super-resolution. *ArXiv*, abs/2103.14006, 2021. 21
- [102] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 3, 7, 16
- [103] Shengyu Zhao, Jianwei Cui, Yilun Sheng, Yue Dong, Xiao Liang, Eric I-Chao Chang, and Yan Xu. Large scale image completion via co-modulated generative adversarial networks. *ArXiv*, abs/2103.10428, 2021. 8
- [104] Bolei Zhou, Àgata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40:1452–1464, 2018. 8, 24
- [105] Yufan Zhou, Ruiyi Zhang, Changyou Chen, Chunyuan Li, Chris Tensmeyer, Tong Yu, Jiuxiang Gu, Jinhui Xu, and Tong Sun. LAFITE: towards language-free training for text-to-image generation. *CoRR*, abs/2111.13792, 2021. 6