# Unleashing Transformers: Parallel Token Prediction with Discrete Absorbing Diffusion for Fast High-Resolution Image Generation from Vector-Quantized Codes

Sam Bond-Taylor[*1], Peter Hessey[*1], Hiroshi Sasaki[1], Toby P. Breckon[1,2], Chris G. Willcocks[1]

Department of {[1]Computer Science | [2]Engineering}, Durham University, Durham, UK

## Abstract

*Whilst diffusion probabilistic models can generate high quality image content, key limitations remain in terms of both generating high-resolution imagery and their associated high computational requirements. Recent Vector-Quantized image models have overcome this limitation of image resolution but are prohibitively slow and unidirectional as they generate tokens via element-wise autoregressive sampling from the prior. By contrast, in this paper we propose a novel discrete diffusion probabilistic model prior which enables parallel prediction of Vector-Quantized tokens by using an unconstrained Transformer architecture as the backbone. During training, tokens are randomly masked in an order-agnostic manner and the Transformer learns to predict the original tokens. This parallelism of Vector-Quantized token prediction in turn facilitates unconditional generation of globally consistent high-resolution and diverse imagery at a fraction of the computational expense. In this manner, we can generate image resolutions exceeding that of the original training set samples whilst additionally provisioning per-image likelihood estimates (in a departure from generative adversarial approaches). Our approach achieves state-of-the-art results in terms of Density (LSUN Bedroom: 1.51; LSUN Churches: 1.12; FFHQ: 1.20) and Coverage (LSUN Bedroom: 0.83; LSUN Churches: 0.73; FFHQ: 0.80), and performs competitively on FID (LSUN Bedroom: 3.64; LSUN Churches: 4.07; FFHQ: 6.11) whilst offering advantages in terms of both computation and reduced training set requirements.*

## 1. Introduction

Artificially generating plausible photo-realistic images, at ever higher resolutions, has long been a goal when designing deep generative models. Recent advancements have yielded direct benefits for fields such as medical image synthesis [19], computer graphics [9, 79], image editing



Figure 1. Our approach uses a discrete diffusion to generate high quality images optionally larger than the training data (bottom).

[45], image-to-image translation [65], and image super-resolution [27].

These methods can in general be divided into five main classes [5], each of which make different trade-offs to scale to high resolutions. Techniques to scale Generative Adversarial Networks (GANs) [20] include progressive growing [34], large batches [8], and regularisation [46, 49]. Variational Autoencoders (VAEs) [43] can be scaled by building complex priors [10, 70, 73] and correcting the learned density [77]. Autoregressive approaches can make independence assumptions [60] or partition spatial dimensions [48]. Normalizing Flows utilise multi-scale architectures [40], while diffusion models can be scaled using SDEs [67] and cascades [27]. Each of these approaches have their own drawbacks, such as unstable training, long sample times, and a lack of global context.

Of particular interest to this work is the popular Trans-

---

former architecture [74] which is able to model long distance relationships using a powerful attention mechanism that can be trained in parallel. By constraining the Transformer architecture to attend a fixed ordering of tokens in a unidirectional manner, they can be used to parameterise an autoregressive model for generative modelling [11, 55]. However, image data does not conform to such a structure and hence this bias limits the representation ability of the Transformer and unnecessarily restricts the sampling process to be both sequential and slow.

Addressing these issues, our main contributions are:

- We propose a novel parallel token prediction approach for generating Vector-Quantized image representations that allows for significantly faster sampling than autoregressive approaches.

- Our approach is able to generate globally consistent images at resolutions exceeding that of the original training data by aggregating multiple context windows, allowing for much larger context regions.

- Our approach demonstrates state-of-the art performance across three benchmark datasets in terms of Density (LSUN Bedroom: 1.51; LSUN Churches: 1.12; FFHQ: 1.20) and Coverage (LSUN Bedroom: 0.83; LSUN Churches: 0.73; FFHQ: 0.80), while also being competitive on FID (LSUN Bedroom: 3.64; LSUN Churches: 4.07; FFHQ: 6.11).

## 2. Prior Work

Extensive work in deep generative modelling [5] and self-supervised learning [16] laid the foundations for this research, which we review here in terms of both existing models (Sections 2.1-2.4) and the Transformer architecture (Section 2.5).

### 2.1. Autoregressive Models

Autoregressive models are a family of powerful generative models capable of directly maximising the likelihood of the data on which they are trained. These models have achieved impressive image generation results in recent years, however, their sequential nature limits them to relatively low dimensional data [12, 33, 62, 64, 71, 72].

The training and inference process for autoregressive models is based on the chain rule of probability. By decomposing inputs into their individual components $\boldsymbol{x} = \{x_1, ..., x_n\}$, an autoregressive model with parameters $\theta$ can generate new latent samples sequentially:

$$p_\theta(\boldsymbol{x}) = p_\theta(x_1, x_2, ..., x_n) = \prod_{i=1}^{n} p_\theta(x_i | x_1, ..., x_{i-1}). \quad (1)$$

Selecting an ordering over inputs is not obvious for many tasks; since the receptive field is limited to previously generated tokens, this can significantly affect sample quality.

### 2.2. Discrete Energy-Based Models

Since the causal nature of autoregressive models limits their representation ability, other approaches with less constrained architectures have begun to outperform them even on likelihood [42]. Energy-based models (EBMs) are an enticing method for representing discrete data as they permit unconstrained architectures with global context. Implicit EBMs define an unnormalised distribution over data that is typically learned through contrastive divergence [15, 25]. Unfortunately, sampling EBMs using Gibbs sampling is impractical for high dimensional discrete data. However, gradients can be incorporated to reduce mixing times [23].

Similar to autoregressive models, masked language models (MLMs) such as BERT [13] model the conditional probability of the data. However, these are trained bidirectionally by randomly masking a subset of tokens from the input sequence, allowing a much richer context than autoregressive approaches. Some attempts have been made to define an implicit energy function using the conditional probabilities [75], however, obtaining true samples leads to very long sample times and we found them to be ineffective at modelling longer sequences during our experiments [21].

### 2.3. Discrete Denoising Diffusion Models

Diffusion models [26, 66] define a Markov chain $q(\boldsymbol{x}_{1:T} | \boldsymbol{x}_0) = \prod_{t=1}^{T} q(\boldsymbol{x}_t | \boldsymbol{x}_{t-1})$ that gradually destroys data $\boldsymbol{x}_0$ by adding noise over a fixed number of steps $T$ so that $\boldsymbol{x}_T$ contains little to no information about $\boldsymbol{x}_0$ and can be easily sampled. The reverse procedure is a generative model that gradually denoises towards the data distribution $p_\theta(\boldsymbol{x}_{0:T}) = p_\theta(\boldsymbol{x}_T) \prod_{t=1}^{T} p_\theta(\boldsymbol{x}_{t-1} | \boldsymbol{x}_t)$, learned by optimising the evidence lower bound (ELBO), with $t^{\text{th}}$ term

$$\mathbb{E}_{q(\boldsymbol{x}_{t+1} | \boldsymbol{x}_0)} \left[ D_{KL}(q(\boldsymbol{x}_t | \boldsymbol{x}_{t+1}, \boldsymbol{x}_0) || p_\theta(\boldsymbol{x}_t | \boldsymbol{x}_{t+1})) \right], \quad (2)$$

where sampling from the reverse process is not required during training. When applied in continuous spaces, distributions are typically parameterised as Normal distributions.

Discrete diffusion models [1, 28, 66] constrain the state space so that $\boldsymbol{x}_t$ is a discrete random variable falling into one of $K$ categories. As such, the forward process can be represented as categorical distributions $q(\boldsymbol{x}_t | \boldsymbol{x}_{t-1}) = \text{Cat}(\boldsymbol{x}_t; \boldsymbol{p} = \boldsymbol{x}_{t-1} \boldsymbol{Q}_t)$ for one-hot $\boldsymbol{x}_{t-1}$ where $\boldsymbol{Q}_t$ is a matrix denoting the probabilities of moving to each successive state. Transition processes include moving states with some low uniform probability [28], moving to nearby states with some low probability based on similarity or distance, and masking out values entirely similar to generative MLMs.

### 2.4. Hybrid Generative Models

Hybrid models combine two or more classes of generative model to balance trade-offs such as slow sampling, poor scaling with dimension, and inadequate modelling flexibility. Many of these approaches are based on Variational
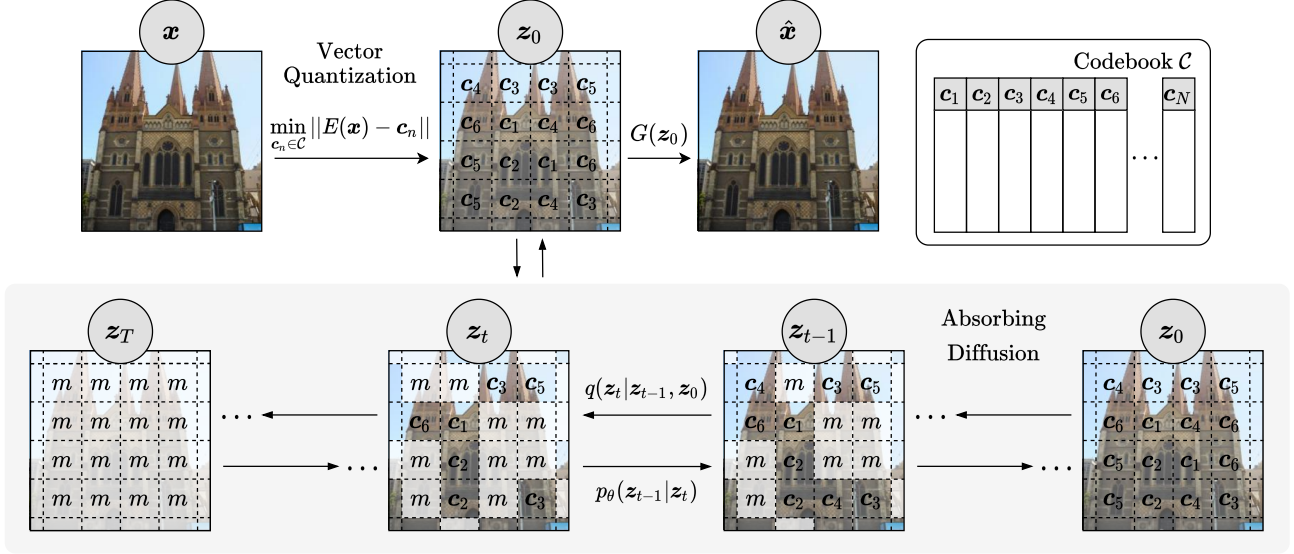
Figure 2. Our approach uses a discrete absorbing diffusion model to represent Vector-Quantized images allowing fast high-resolution image generation. Specifically, after compressing images to an information-rich discrete space, elements are randomly masked and an unconstrained Transformer is trained to predict the original data, using global context to ensure samples are consistent and high quality.

Autoencoders (VAEs) [43, 61] which have fast run-times and scale well to high resolutions, but struggle with sample quality. For example, a VAE's approximate posterior and/or prior complexity can be increased by applying a second generative model such as a Normalizing Flow [22,41,69,70] or EBM [54] in latent space. Alternatively, a second model can be used to correct samples [77]. Of particular interest to this work are Vector-Quantized image models, which follow a 2-stage training scheme where a convolutional autoencoder extracts high level features to an information rich discrete latent space and a powerful autoregressive density estimator learns the prior over these latents [18, 59, 73].

## 2.5. Transformers

Transformers [74] have made a huge impact across many fields of deep learning [24] due to their power and flexibility. They are based on the concept of self-attention, a function which allows interactions with strong gradients between all inputs, irrespective of their spatial relationships. This procedure (Eqn. 3) encodes inputs as key-value pairs, where values $V$ represent embedded inputs and keys $K$ act as an indexing method, subsequently, a set of queries $Q$ are used to select which values to observe:

$$\text{Attn}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V. \quad (3)$$

While this allows long distance dependencies to be learned, complexity increases with sequence length quadratically, making scaling to high dimensional inputs difficult. Approaches to mitigate this include independence assumptions [60], sparsity [11], and kernel approximations [38].

## 3. Method

In this section, we formalise our proposed 2-stage approach for generating high-resolution images using a discrete diffusion model to represent Vector-Quantized image representations; this is visualised in Fig. 2. We hypothesise that by removing the autoregressive constraint, allowing bidirectional context when generating samples, not only will it be possible to speed up the sampling process, but an improved feature representation will be learned, enabling higher quality image generation.

### 3.1. Learning Codes

In the first stage of our approach, a Vector-Quantized image model compresses high-resolution images to a highly compressed form, taking advantage of an information rich codebook [73]. A convolutional encoder downsamples images $x$ to a smaller spatial resolution, $E(x) = \{e_1, e_2, ..., e_L\} \in \mathbb{R}^{L \times D}$. A simple quantisation approach is to use the `argmax` operation which maps continuous encodings to their closest elements in a finite codebook of vectors [73]. Specifically, for a codebook $\mathcal{C} \in \mathbb{R}^{K \times D}$, where $K$ is the number of discrete codes in the codebook and $D$ is the dimension of each code, each $e_i$ is mapped via a nearest-neighbour lookup onto a discrete codebook value, $c_j \in \mathcal{C}$:

$$z_q = \{q_1, q_2, ..., q_L\}, \text{ where } q_i = \min_{c_j \in \mathcal{C}} ||e_i - c_j||. \quad (4)$$

As this operation is non-differentiable, the straight-through gradient estimator [3] is used to copy the gradients from the decoder inputs onto the encoder outputs resulting in bi-

ased gradients. Subsequently, the quantized latents are fed through a decoder network $\hat{\boldsymbol{x}} = G(\boldsymbol{z}_q)$ to reconstruct the input based on a perceptual reconstruction loss [18, 80]; this process is trained by minimising the loss $\mathcal{L}_{\text{VQ}}$,

$$\mathcal{L}_{\text{VQ}} = \mathcal{L}_{\text{rec}} + ||\text{sg}[E(\boldsymbol{x})] - \boldsymbol{z}_q||_2^2 + \beta ||\text{sg}[\boldsymbol{z}_q] - E(\boldsymbol{x})||_2^2. \quad (5)$$

The $\text{argmax}$ approach can result in codebook collapse, where some codes are never used; while other quantisation methods can reduce this [14, 32, 47, 58], we found $\text{argmax}$ quantisation to yield the highest reconstruction quality.

### 3.2. Sampling Globally Coherent Latents

To allow sampling, a discrete generative model is trained on the latents obtained from the Vector-Quantized image model. The highly compressed form allows this second stage to function much more efficiently. Once the training data is encoded as discrete, integer-valued latents $\boldsymbol{z} \in \mathbb{Z}^D$, a discrete diffusion model can be used to learn the distribution over these latents. Due to the effectiveness of BERT-style models [13] for representation learning, we use the absorbing state diffusion [1] which similarly learns to denoise randomly masked data. Specifically, in each forward time step $t$, values are either kept the same or masked out entirely with probability $\frac{1}{t}$ and the reverse process gradually unveils these masks. Rather than directly approximating $p_\theta(\boldsymbol{z}_{t-1}|\boldsymbol{z}_t)$, we predict $p_\theta(\boldsymbol{z}_0|\boldsymbol{z}_t)$, reducing the training stochasticity [26]. The variational bound reduces to

$$\mathbb{E}_{q(\boldsymbol{z}_0)} \left[ \sum_{t=1}^{T} \frac{1}{t} \mathbb{E}_{q(\boldsymbol{z}_t|\boldsymbol{z}_0)} \left[ \sum_{[\boldsymbol{z}_t]_i = m} \log p_\theta([\boldsymbol{z}_0]_i | \boldsymbol{z}_t) \right] \right]. \quad (6)$$

In practice, continuous diffusion models are trained to estimate the noise rather than directly predict the denoised data; this reparameterisation allows the loss to be easily minimised at time steps close to $T$. Unfortunately, no relevant reparameterisation currently exists for discrete distributions [28]. Rather than directly maximising the ELBO, we reweight the ELBO to mimic the reparameterisation,

$$\mathbb{E}_{q(\boldsymbol{z}_0)} \left[ \sum_{t=1}^{T} \frac{T-t+1}{T} \mathbb{E}_{q(\boldsymbol{z}_t|\boldsymbol{z}_0)} \left[ \sum_{[\boldsymbol{z}_t]_i = m} \log p_\theta([\boldsymbol{z}_0]_i | \boldsymbol{z}_t) \right] \right], \quad (7)$$

where components of the loss at time steps close to $T$ are weighted less than earlier steps. This is closely related to the loss obtained by assuming the posterior does not have access to $\boldsymbol{x}_t$, i.e. if the $t-1^{\text{th}}$ loss term is $D_{KL}(q(\boldsymbol{x}_{t-1}|\boldsymbol{x}_0)||p_\theta(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t))$. Since we directly predict $\boldsymbol{z}_0$ and not $\boldsymbol{z}_t$ this assumption does not harm the training. Experimentally we find that this reweighting achieves lower validation ELBO than when directly maximising the ELBO.

Esser et al. [18] demonstrated that in the autoregressive case, Transformers [74] are better suited for modelling Vector-Quantized images than convolutional architectures

due to the importance of long-distance relationships in this compressed form. As such, we utilise transformers to model the prior distribution, but without the architectural restrictions imposed by autoregressive approaches.

### 3.3. Generating High-Resolution Images

Using convolutions to build Vector-Quantized image models encourages latents to be highly spatially correlated with generated images. It is therefore possible to construct essentially arbitrarily sized images by generating latents with the required shape. We propose an approach that allows globally consistent images substantially larger than those in the training data to be generated.

First, a large $a$ by $b$ array of mask tokens, $\bar{\boldsymbol{z}}_T = m^{a \times b}$, is initialised that corresponds to the size of image we wish to generate. In order to capture the maximum context when approximating $\bar{\boldsymbol{z}}_0$ we apply the denoising network to all subsets of $\bar{\boldsymbol{z}}_t$ with the same spatial size as the usual inputs of the network, aggregating estimates at each location. Specifically, using $c_j(\bar{\boldsymbol{z}}_t)$ to represent local subsets, we approximate the denoising distribution as a mixture,

$$p([\bar{\boldsymbol{z}}_0]_i | \bar{\boldsymbol{z}}_t) \approx \frac{1}{Z} \sum_j p([\bar{\boldsymbol{z}}_0]_i | c_j(\bar{\boldsymbol{z}}_t)), \quad (8)$$

where the sum is over subsets $c_j$ that contain the $i^{th}$ latent. For extremely large images, this can require a very large number of function evaluations, however, the sum can be approximated by striding over latents with a step $> 1$ or by randomly selecting positions.

### 3.4. Improving Code Representations

There are various options to obtain high-quality image representations including using large numbers of latents and codes [58] or building a hierarchy of latent variables [59]. We use the adversarial framework proposed by Esser et al. [18] to achieve higher compression rates with high-quality codes using only a single GPU, without tying our approach to the characteristics typically associated with generative adversarial models. Additionally, we apply differentiable augmentations $T$, such as translations and colour jitter, to all discriminator inputs; this has proven to be effective at improving sample quality across methods [33, 81]. The overall loss $\mathcal{L}$ is a linear combination of $\mathcal{L}_{\text{VQ}}$, the Vector-Quantized loss, and $\mathcal{L}_{\text{G}}$ which uses a discriminator $D$ to assess realism based on an adaptive weight $\lambda$. On some datasets, $\lambda$ can grow to extremely large values hindering training. We find simply clamping $\lambda$ at a maximum value $\lambda_{\text{max}} = 1$ an effective solution that stabilises training,

$$\mathcal{L} = \min_{E,G,\mathcal{C}} \max_{D} \mathbb{E}_{\boldsymbol{x} \sim p_d} \left[ \mathcal{L}_{\text{VQ}} + \lambda \mathcal{L}_{\text{G}} \right], \quad (9a)$$

$$\mathcal{L}_{\text{G}} = \log D(T(\boldsymbol{x})) + \log(1 - D(T(\hat{\boldsymbol{x}}))), \quad (9b)$$

$$\lambda = \min \left( \frac{\nabla_{G_L}[\mathcal{L}_{\text{rec}}]}{\nabla_{G_L}[\mathcal{L}_{\text{G}}] + \delta}, \lambda_{\text{max}} \right). \quad (9c)$$

| Model | P ↑ | R ↑ | D ↑ | C ↑ |
|---|---|---|---|---|
| **Churches** | | | | |
| DCT [51] | 0.60 | **0.48** | - | - |
| TT [18] | 0.67 | 0.29 | 1.08 | 0.60 |
| PGGAN [34] | 0.61 | 0.38 | 0.83 | 0.63 |
| StyleGAN2 [37] | 0.60 | 0.43 | 0.83 | 0.68 |
| **Ours** ($t = 1.0$) | 0.70 | 0.42 | **1.12** | 0.73 |
| **Ours** ($t = 0.9$) | **0.71** | 0.45 | 1.07 | **0.74** |
| **FFHQ** | | | | |
| VDVAE [10] | 0.59 | 0.20 | 0.80 | 0.50 |
| TT [18] | 0.64 | 0.29 | 0.89 | 0.59 |
| StyleGAN2 [37] | 0.69 | 0.40 | 1.12 | 0.80 |
| **Ours** ($t = 1.0$) | 0.69 | 0.48 | 1.06 | 0.77 |
| **Ours** ($t = 0.9$) | **0.73** | **0.48** | **1.20** | **0.80** |
| **Bedroom** | | | | |
| DCT [51] | 0.44 | **0.56** | - | - |
| TT [18] | 0.61 | 0.33 | 1.15 | 0.75 |
| PGGAN [34] | 0.43 | 0.40 | 0.70 | 0.64 |
| StyleGAN [36] | 0.55 | 0.48 | 0.96 | 0.80 |
| **Ours** ($t = 1.0$) | 0.64 | 0.38 | 1.27 | 0.81 |
| **Ours** ($t = 0.9$) | **0.67** | 0.38 | **1.51** | **0.83** |

Table 1. Precision, Recall, Density, and Coverage for approaches trained on FFHQ, LSUN Bedroom, and LSUN Churches.

| Method | Params | Bed | Church | FFHQ |
|---|---|---|---|---|
| DDPM [26] | 114M | 6.36 | 7.89 | - |
| DCT [51] | 448M | 6.40 | 7.56 | - |
| VDVAE [10] | 115M | - | - | 28.5 |
| TT [17, 18] | 600M | 6.35 | 7.81 | 9.6 |
| ImageBART [17] | 2104M | 5.51 | 7.32 | 9.57 |
| PGGAN [34] | 47M | 8.34 | 6.42 | - |
| StyleGAN2 [37] | 60M | 2.35 | 3.86 | 3.8 |
| **Ours** ($t = 1.0$) | 145M | 5.07 | 5.58 | 7.12 |
| **Ours** ($t = 0.9$) | 145M | 3.64 | 4.07 | 6.11 |

Table 2. FID for various approaches on FFHQ, LSUN Bedroom, and LSUN Churches. Lower FID signifies higher quality samples.

# 4. Evaluation

We evaluate our approach on three high-resolution 256x256 datasets: LSUN Bedroom, LSUN Churches [78], and FFHQ [36]. Sec. 4.1 evaluates the quality of samples from our proposed model. Sec. 4.2 demonstrates the representation abilities of absorbing diffusion models applied to the learned discrete latent spaces, including how sampling can be sped up, improvements over equivalent autoregressive models, and the effect of our reweighted ELBO. Finally, Sec. 4.4 evaluates our Vector-Quantized image model.

In all experiments, our absorbing diffusion model parameterised with an 80M parameter Transformer Encoder



Figure 3. Samples from our models trained on 256x256 datasets: LSUN Churches, FFHQ, and LSUN Bedroom.

[74] is applied to $16 \times 16$ latents discretised to a codebook with 1024 entries and optimised using the Adam optimiser [39]. While, as noted by Esser et al. [18], a GPT2-medium [57] architecture (307M parameters) fits onto a GPU with 12GB of VRAM, in practice this requires the use of small batch sizes and learning rates making training in reasonable times impractical. More details can be found in Appendix A. Source code for the models used and experiments performed in this paper is available here.

## 4.1. Sample Quality

In this section we evaluate samples from our model quantitatively and qualitatively. In comparison to other multi-step methods, our approach allows sampling in the fewest steps. Samples from our model can be found in Fig. 3 which are high quality and diverse. More samples can be found in Appendix D.

**PRDC** In Tab. 1 we evaluate our approach against a variety of other models in terms of Precision, Recall, Density, and Coverage (PRDC) [44, 50, 63], metrics that quantify the overlap between the data and sample distributions. Due to limited computing resources, we are unable to provide density and coverage scores for DCT [51] and PRDC scores for StyleGAN2 on LSUN Bedroom since training on a standard GPU would take more than 30 days per experiment, signif-

Figure 4. Our method allows unconditional images larger than those seen during training to be generated by applying the denoising network to all subsets of the image, aggregating probabilities to encourage global continuity.

icantly more than the 10 days required to train our models. On the LSUN datasets our approach achieves the highest Precision, Density, and Coverage; indicating that the data and sample manifolds have the most overlap. On FFHQ our approach achieves the highest Precision and Recall. In general, when generative models are sampled with lower temperatures to achieve lower FID, this leads to trade-off between precision and recall [37, 59]; since we also calculate FID with a lower temperature, we evaluate the effect of this on PRDC. In all but one case sampling with temperature leads to improved scores, indicating that our approach fits more accurately.

**FID** In Tab. 2 we calculate the Fréchet Inception Distance (FID) of samples from our models using torch-fidelity [53]. Despite using a fraction of the number of parameters compared to other Vector-Quantized image models, our approach achieves substantially lower FID scores.

**Higher Resolution** Fig. 4 contains samples of various spatial sizes using the approach described in Sec. 3.3. Here, an absorbing diffusion model is trained on $16 \times 16$ latents then samples are generated at larger sizes (up to $768 \times 256$) using a temperature value of 0.8. Even at the larger scales we observe high-quality, diverse, and consistent imagery.

#### 4.1.1 Limitations of FID Metric

While FID has been found to correlate well with image quality, it unrealistically approximates the data distribution as Gaussian in embedding space and is insensitive to the global structure of the data distribution [68]. For likelihood models, calculating NLL on a test set is possible instead but likelihood has been shown to not correlate well with quality; fine tuning our approach to model pixels as Gaussians gives 2.72BPD on 5-bit FFHQ. Alternative approaches that address these issues have been developed [6] such as PPL [36], which assesses sample consistency through latent interpolations; IMD [68], which uses all moments to compare data manifolds making it sensitive to global structure; and MTD [2], which compares manifolds in image space.

In this work, we compare approaches using Precision and Recall [63] approaches which, unlike FID, evaluate sample quality and diversity separately and have been used in similar recent work assessing high-resolution image generation [30, 37, 51, 59]. Precision is the expected likelihood of fake samples lying on the data manifold and recall vice versa. These metrics are computed by approximating the data and sample manifolds as hyper-spheres around data and sample points respectively [44]. Density and Coverage are modifications to Precision and Recall respectively that address manifold overestimation [50].

### 4.2. Absorbing Diffusion

In this section we analyse the usage of absorbing diffusion for high-resolution image generation, determining how many sampling steps are required to obtain high-quality samples and ablating the components of our approach.

#### 4.2.1 Sampling Speed

Our approach applies a diffusion process to a highly compressed image representation, meaning it is already $18\times$ faster to sample from than DDPM (ours: 3.8s, DDPM: 70s per image on a NVIDIA RTX 2080 Ti). However, since the absorbing diffusion model is trained to approximate $p(\boldsymbol{z}_0|\boldsymbol{z}_t)$ it is possible to speed the sampling process up further by skipping arbitrary numbers of time steps, unmasking multiple latents at once. In Tab. 3 we explore how sample quality is affected using a simple step skipping scheme: evenly skipping a constant number of steps so that the to-

| Steps | 50 | 100 | 150 | 200 | 256 |
|---------|------|------|------|------|------|
| Churches | 6.86 | 6.09 | 5.81 | 5.68 | 5.58 |
| Bedroom | 6.85 | 5.83 | 5.53 | 5.32 | 5.42 |
| FFHQ | 9.60 | 7.90 | 7.53 | 7.52 | 7.12 |

Table 3. FID for different number of sampling steps on LSUN Churches, Bedroom and FFHQ. Diffusion steps are evenly spaced.

| Method | Churches | | FFHQ | |
|---|---|---|---|---|
| | FID ↓ | NLL ↓ | FID ↓ | NLL ↓ |
| Autoregressive | 5.93 | 6.24 | 8.15 | 6.18 |
| Absorbing DDPM | **5.58** | **6.01** | **7.12** | **5.96** |

Table 4. FID and validation NLL (in BPD) for different methods to approximate discrete latents using the same Transformer architecture on LSUN Churches and FFHQ.
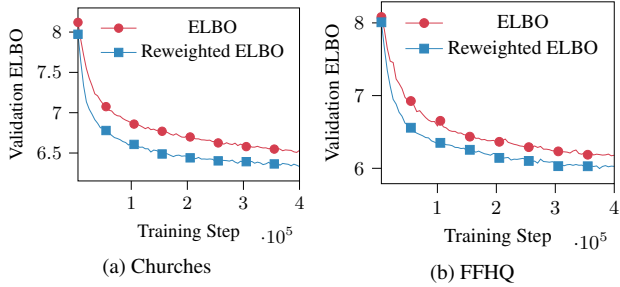


Figure 5. Comparison of proposed losses on (a) LSUN Churches and (b) FFHQ. Models trained with our reweighted ELBO achieve lower validation ELBO than models trained directly on the ELBO.

tal number of steps meets some fixed computational budget. As expected, FID increases with fewer sampling steps. However, the increase in FID is minor relative to the improvement in sampling speed: our approach achieves similar FID to the equivalent autoregressive model using half the number of steps. With 50 sampling steps, our approach is $88\times$ faster than DDPMs. Using a more sophisticated step selection scheme such as dynamic programming [76], FID could potentially be reduced further.

### 4.2.2 Autoregressive vs Absorbing DDPM

Tab. 4 compares the representation ability of our absorbing diffusion model with an autoregressive model, both utilising exactly the same Transformer architecture, but with the Transformer unconstrained in the diffusion case. On both datasets tested, the diffusion achieves lower FID and NLL than the autoregressive model despite being trained on a harder task with the same number of parameters; the additional regularisation prevents overfitting which is prevalent with autoregressive models [17, 33]. Since the number of sampling steps is the same as the number of data dimensions, samples from the diffusion are effectively autoregressive over random orderings, indicating that the learned distributions better approximate the data distribution.

### 4.2.3 Reweighted ELBO

In Sec. 3.2 we proposed using a reweighted ELBO when training the diffusion model that focuses gradients on the central training steps, balancing feasability with ease of
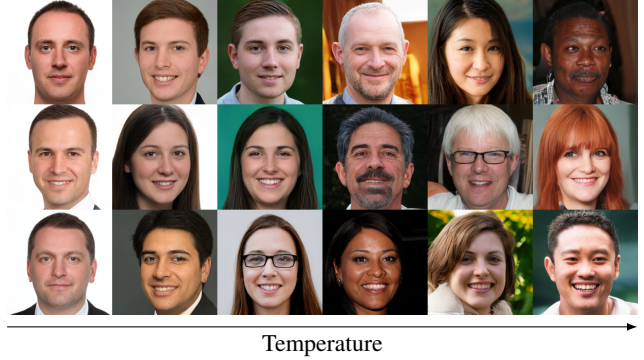


Figure 6. Impact of sampling temperature on diversity. For small temperature changes it is less obvious how bias has changed.

learning. We evaluate this in Fig. 5 by comparing validation ELBO during training for models trained directly on ELBO and our re-weighting. The re-weighted ELBO converges to a lower validation ELBO sooner, demonstrating that our reweighting is valid and simplifies optimisation.

### 4.3. Sample Diversity

To improve sample quality, many generative models are sampled using a reduced temperature or by truncating Normal distributions. This is problematic, as these sampling methods will amplify any biases in the dataset. We visualise the impact of temperature on sampling from a model trained on FFHQ in Fig. 6. For very low temperatures the bias the obvious: samples are mostly front-facing white men with brown hair on solid white/black backgrounds. Exactly how the bias has changed for more subtle temperature changes is less clear, which is problematic. Practitioners should be aware of this effect and it emphasises the importance of dataset balancing.

### 4.4. Reconstruction Quality

In Tab. 5 we evaluate the effect of differentiable augmentations (DiffAug) [81] and adaptive weight limiting on Vector-Quantized image modelling. While applying each technique individually can lead to worse FID due to imbalance, when both techniques are applied, we found that FID improved across all datasets tested.

| Modifications | Churches | FFHQ |
|---|---|---|
| Default | 5.25 | 3.37 |
| $\lambda_{\max} = 1$ | 8.67 | 4.72 |
| DiffAug | 5.16 | 6.57 |
| Both | **2.70** | **3.12** |

Table 5. Effect of DiffAug and adaptive weight limiting on reconstruction FID. Results were calculated on Vector-Quantized image models trained on LSUN Churches and FFHQ for 500k steps.
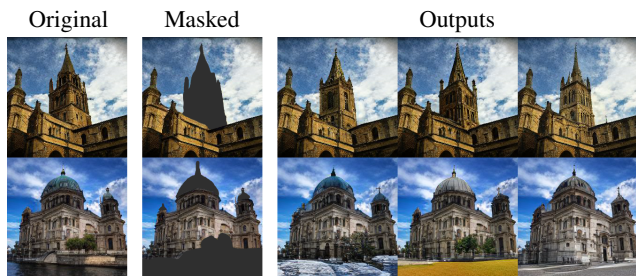
Original     Masked         Outputs

Figure 7. Our approach allows local image editing by targeting regions to be changed (highlighted in grey). Here alternative towers are generated and water is replaced with different foregrounds.

### 4.5. Image Editing

An additional advantage of using a bidirectional diffusion model to model the latent space is that image inpainting is possible. Since autoregressive models are conditioned only on the upper left region of the image, they are unable to edit internal masked image regions in a consistent manner. Diffusion models, on the other hand, allow masked regions to be placed at arbitrary locations. After a region has been highlighted, we mask corresponding latents, identify the starting time step by counting the number of masked latents, then continue the denoising process from that point. Examples of this process can be found in Fig. 7.

### 4.6. Limitations

In our experiments we only tested our approach on 256×256 datasets; directly scaling to higher resolutions would require more GPU resources. However, future work using more efficient Transformer architectures [31] may alleviate this. Our method outperforms all approaches tested on FID except StyleGAN2 [37]; we find that the primary bottleneck is the Vector-Quantized image model, therefore more research is necessary to improve these discrete representations. Whilst our approach is trained for significantly less time than other approaches such as StyleGAN2, the stochastic training procedure means that more training steps are required compared to autoregressive approaches. Although when generating extra-large images the large context window made possible by the diffusion model encourages consistency, a reduced temperature is also required, reducing diversity.

### 5. Discussion

While other classes of discrete generative model exist, they are less suitable for Vector-Quantized image modelling than discrete diffusion models: VAEs introduce prior assumptions about the latent space that can be limiting, in particular, continuous spaces may not be appropriate when modelling discrete data [7]; GAN training requires sampling from the generator meaning that gradients must be back-

propagated through a discretistion procedure [52]; Discrete normalising flows require the use of invertible functions which significantly restrict the function space [4, 29].

Concurrently developed with this work was ImageBART [17] which also uses a diffusion model to learn the prior of a Vector-Quantized image model. However, our approaches substantially differ: ImageBART uses a multinomial diffusion process with separate autoregressive Transformers trained to approximate each diffusion step, leading to slower inference and substantially more parameters; our approach optimises all diffusion steps simultaneously with a single, non-autoregressive Transformer.

There are numerous potential avenues to explore to further improve sample quality. For instance, when learning discrete image representations, implicit networks (which are invariant to translation and rotation) [35] or other more powerful generative models could be used. Alternatively, different discrete diffusion methods could be used that impose relationships between codes based on their continuous embeddings. Finally, by conditioning on both text and discrete image representations, absorbing diffusion models could allow text-to-image generation and image captioning to be accomplished using a single model with faster runtime than independent approaches [56, 58].

### 5.1. Social Impact

While deep generative models have various positive applications such as text-to-speech, drug design, and generating examples of rare medical conditions, there can be negative consequences associated with their development:

- As with all generative models, our approach can be used for malicious applications such as deep-fakes.
- Samples reflect the biases present in the training data which can lead to unintended consequences.
- Training generative models such as ours consumes significant quantities of energy affecting the environment. The fast sampling permitted by our approach, however, reduces this at test-time compared to similar methods.

### 6. Conclusion

In this work we proposed a discrete diffusion probabilistic model prior capable of predicting Vector-Quantized image representations in parallel, overcoming the high sampling times, unidirectional nature and overfitting challenges associated with autoregressive priors. Our approach makes no assumptions about the inherent ordering of latents by utilising an unconstrained Transformer architecture. Experimental results demonstrate the ability of our approach to generate diverse, high-quality images, optionally at resolutions exceeding the training samples. Additional work is needed to reduce training times and to efficiently scale our approach to even higher resolutions.

# References

[1] Jacob Austin, Daniel Johnson, Jonathan Ho, Danny Tarlow, and Rianne van den Berg. Structured Denoising Diffusion Models in Discrete State-Spaces. *arXiv preprint arXiv:2107.03006*, 2021. 2, 4

[2] Serguei Barannikov, Ilya Trofimov, Grigorii Sotnikov, Ekaterina Trimbach, Alexander Korotin, Alexander Filippov, and Evgeny Burnaev. Manifold Topology Divergence: a Framework for Comparing Data Manifolds. *arXiv preprint arXiv:2106.04024*, 2021. 6

[3] Yoshua Bengio. Estimating or propagating gradients through stochastic neurons, 2013. 3

[4] Rianne van den Berg, Alexey A Gritsenko, Mostafa Dehghani, Casper Kaae Sønderby, and Tim Salimans. IDF++: Analyzing and Improving Integer Discrete Flows for Lossless Compression. In *International Conference on Learning Representations*, 2021. 8

[5] Sam Bond-Taylor, Adam Leach, Yang Long, and Chris G Willcocks. Deep Generative Modelling: A Comparative Review of VAEs, GANs, Normalizing Flows, Energy-Based and Autoregressive Models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 1, 2

[6] Ali Borji. Pros and Cons of GAN Evaluation Measures: New Developments. *arXiv preprint arXiv:2103.09396*, 2021. 6

[7] Samuel R. Bowman, Luke Vilnis, Oriol Vinyals, Andrew M. Dai, Rafal Jozefowicz, and Samy Bengio. Generating Sentences from a Continuous Space. *arXiv:1511.06349*, 2016. 8

[8] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large Scale GAN Training for High Fidelity Natural Image Synthesis. In *International Conference on Learning Representations*, 2019. 1

[9] Xuelin Chen, Daniel Cohen-Or, Baoquan Chen, and Niloy J Mitra. Towards a Neural Graphics Pipeline for Controllable Image Generation. In *Computer Graphics Forum*, volume 40, pages 127–140. Wiley Online Library, 2021. 1

[10] Rewon Child. Very Deep VAEs Generalize Autoregressive Models and Can Outperform Them on Images. In *International Conference on Learning Representations*, 2021. 1, 5

[11] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating Long Sequences with Sparse Transformers. *arXiv preprint arXiv:1904.10509*, 2019. 2, 3

[12] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating Long Sequences with Sparse Transformers. *arXiv:1904.10509*, 2019. 2

[13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-Training of Deep Bidirectional Transformers for Language Understanding. In *NAACL-HLT*, 2019. 2, 4

[14] Sander Dieleman, Aäron van den Oord, and Karen Simonyan. The Challenge of Realistic Music Generation: Modelling Raw Audio at Scale. In *Advances in Neural Information Processing Systems*, volume 31, 2018. 4

[15] Yilun Du and Igor Mordatch. Implicit Generation and Generalization in Energy-Based Models. In *Advances in Neural Information Processing Systems*, volume 33, 2019. 2

[16] Linus Ericsson, Henry Gouk, and Timothy M Hospedales. How Well Do Self-Supervised Models Transfer? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5414–5423, 2021. 2

[17] Patrick Esser, Robin Rombach, Andreas Blattmann, and Björn Ommer. Imagebart: Bidirectional Context with Multinomial Diffusion for Autoregressive Image Synthesis. *arXiv preprint arXiv:2108.08827*, 2021. 5, 7, 8

[18] Patrick Esser, Robin Rombach, and Björn Ommer. Taming Transformers for High-Resolution Image Synthesis. *arXiv:2012.09841*, 2021. 3, 4, 5, 12

[19] Lukas Fetty, Mikael Bylund, Peter Kuess, Gerd Heilemann, Tufve Nyholm, Dietmar Georg, and Tommy Löfstedt. Latent Space Manipulation for High-Resolution Medical Image Synthesis via the StyleGAN. *Zeitschrift für Medizinische Physik*, 30(4):305–314, 2020. 1

[20] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative Adversarial Nets. In *Advances in Neural Information Processing Systems*, volume 27, 2014. 1

[21] Kartik Goyal, Chris Dyer, and Taylor Berg-Kirkpatrick. Exposing the Implicit Energy Networks behind Masked Language Models via Metropolis–Hastings. *arXiv preprint arXiv:2106.02736*, 2021. 2

[22] Will Grathwohl, Ricky T. Q. Chen, Jesse Bettencourt, Ilya Sutskever, and David Duvenaud. FFJORD: Free-form Continuous Dynamics for Scalable Reversible Generative Models. In *International Conference on Learning Representations*, 2019. 3

[23] Will Grathwohl, Kevin Swersky, Milad Hashemi, David Duvenaud, and Chris J Maddison. Oops I Took A Gradient: Scalable Sampling for Discrete Distributions. In *International Conference on Machine Learning*, 2021. 2

[24] Kai Han, Yunhe Wang, Hanting Chen, Xinghao Chen, Jianyuan Guo, Zhenhua Liu, Yehui Tang, An Xiao, Chunjing Xu, Yixing Xu, et al. A Survey on Visual Transformer. *arXiv preprint arXiv:2012.12556*, 2020. 3

[25] Geoffrey E. Hinton. Training Products of Experts by Minimizing Contrastive Divergence. *Neural Computation*, 14(8):1771–1800, 2002. 2

[26] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models. In *Advances in Neural Information Processing Systems*, volume 33, 2020. 2, 4, 5

[27] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded Diffusion Models for High Fidelity Image Generation. *arXiv preprint arXiv:2106.15282*, 2021. 1

[28] Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax Flows and Multinomial Diffusion: Towards Non-Autoregressive Language Models. *arXiv preprint arXiv:2102.05379*, 2021. 2, 4

[29] Emiel Hoogeboom, Jorn Peters, Rianne van den Berg, and Max Welling. Integer Discrete Flows and Lossless Compression. In *Advances in Neural Information Processing Systems*, volume 32, 2019. 8

[30] Drew A Hudson and C. Lawrence Zitnick. Generative Adversarial Transformers. *Proceedings of the 38th International Conference on Machine Learning, ICML*, 2021. 6

[31] Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David Ding, Skanda Koppula, Daniel Zoran, Andrew Brock, Evan Shelhamer, et al. Perceiver IO: A General Architecture for Structured Inputs & Outputs. *arXiv preprint arXiv:2107.14795*, 2021. 8

[32] Eric Jang, Shixiang Gu, and Ben Poole. Categorical Reparameterization with Gumbel-Softmax. In *International Conference on Learning Representations*, 2017. 4

[33] Heewoo Jun, Rewon Child, Mark Chen, John Schulman, Aditya Ramesh, Alec Radford, and Ilya Sutskever. Distribution Augmentation for Generative Modeling. In *ICML*, 2020. 2, 4, 7

[34] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive Growing of GANs for Improved Quality, Stability, and Variation. In *International Conference on Learning Representations*, 2018. 1, 5

[35] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-Free Generative Adversarial Networks. *arXiv preprint arXiv:2106.12423*, 2021. 8

[36] Tero Karras, Samuli Laine, and Timo Aila. A Style-Based Generator Architecture for Generative Adversarial Networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019. 5, 6

[37] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and Improving the Image Quality of StyleGAN. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 5, 6, 8

[38] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention. In *International Conference on Machine Learning*, 2020. 3

[39] Diederik P Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5, 12

[40] Durk P Kingma and Prafulla Dhariwal. Glow: Generative Flow with Invertible 1x1 Convolutions. In *Advances in Neural Information Processing Systems*, volume 31, 2018. 1

[41] Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improved Variational Inference with Inverse Autoregressive Flow. In *Advances in Neural Information Processing Systems*, volume 29, 2016. 3

[42] Diederik P Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational Diffusion Models. *arXiv preprint arXiv:2107.00630*, 2021. 2

[43] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *International Conference on Learning Representations*, 2014. 1, 3

[44] Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved Precision and Recall Metric for Assessing Generative Models. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32, 2019. 5, 6

[45] Ji Lin, Richard Zhang, Frieder Ganz, Song Han, and Jun-Yan Zhu. Anycost GANs for Interactive Image Synthesis and Editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14986–14996, 2021. 1

[46] Bingchen Liu, Yizhe Zhu, Kunpeng Song, and Ahmed Elgammal. Towards Faster and Stabilized GAN Training for High-fidelity Few-shot Image Synthesis. In *International Conference on Learning Representations*, 2021. 1

[47] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The Concrete Distribution: A Continuous Relaxation of Discrete Random Variables. In *International Conference on Learning Representations*, 2017. 4

[48] Jacob Menick and Nal Kalchbrenner. Generating High Fidelity Images with Subscale Pixel Networks and Multidimensional Upscaling. In *International Conference on Learning Representations*, 2019. 1

[49] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral Normalization for Generative Adversarial Networks. In *International Conference on Learning Representations*, 2018. 1

[50] Muhammad Ferjad Naeem, Seong Joon Oh, Youngjung Uh, Yunjey Choi, and Jaejun Yoo. Reliable Fidelity and Diversity Metrics for Generative Models. In *International Conference on Machine Learning*, pages 7176–7185, 2020. 5, 6, 12

[51] Charlie Nash, Jacob Menick, Sander Dieleman, and Peter W Battaglia. Generating Images with Sparse Representations. *arXiv preprint arXiv:2103.03841*, 2021. 5, 6, 12

[52] Weili Nie, Nina Narodytska, and Ankit Patel. RelGAN: Relational Generative Adversarial Networks for Text Generation. In *International Conference on Learning Representations*, 2019. 8

[53] 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. 6

[54] Bo Pang, Tian Han, Erik Nijkamp, Song-Chun Zhu, and Ying Nian Wu. Learning Latent Space Energy-Based Prior Model. In *Advances in Neural Information Processing Systems*, volume 34, 2020. 3

[55] Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. Image Transformer. In *ICML*, 2018. 2

[56] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning Transferable Visual Models From Natural Language Supervision. *arXiv preprint arXiv:2103.00020*, 2021. 8

[57] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019. 5, 12

[58] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-Shot Text-to-Image Generation. *arXiv preprint arXiv:2102.12092*, 2021. 4, 8

[59] Ali Razavi, Aaron van den Oord, and Oriol Vinyals. Generating Diverse High-Fidelity Images with VQ-VAE-2. *NeurIPS 32*, 2019. 3, 4, 6

[60] Scott Reed, Aäron Oord, Nal Kalchbrenner, Sergio Gómez Colmenarejo, Ziyu Wang, Yutian Chen, Dan Belov, and Nando Freitas. Parallel Multiscale Autoregressive Density Estimation. In *International Conference on Machine Learning*, 2017. 1, 3

[61] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic Backpropagation and Approximate Inference in Deep Generative Models. In *International Conference on Machine Learning*, 2014. 3

[62] Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. Efficient Content-Based Sparse Attention with Routing Transformers. *Transactions of the Association for Computational Linguistics*, 9:53–68, 2021. 2

[63] Mehdi SM Sajjadi, Olivier Bachem, Mario Lucic, Olivier Bousquet, and Sylvain Gelly. Assessing Generative Models via Precision and Recall. In *Advances in Neural Information Processing Systems*, volume 31, 2018. 5, 6

[64] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P. Kingma. PixelCNN++: Improving the PixelCNN with Discretized Logistic Mixture Likelihood and Other Modifications. *ICLR*, 2017. 2

[65] Hiroshi Sasaki, Chris G Willcocks, and Toby P Breckon. UNIT-DDPM: UNpaired Image Translation with Denoising Diffusion Probabilistic Models. *arXiv preprint arXiv:2104.05358*, 2021. 1

[66] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep Unsupervised Learning using Nonequilibrium Thermodynamics. In *International Conference on Machine Learning*, 2015. 2

[67] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-Based Generative Modeling through Stochastic Differential Equations. In *International Conference on Learning Representations*, 2021. 1

[68] Anton Tsitsulin, Marina Munkhoeva, Davide Mottin, Panagiotis Karras, Alex Bronstein, Ivan Oseledets, and Emmanuel Müller. The Shape of Data: Intrinsic Distance for Data Distributions. In *International Conference on Learning Representations*, 2020. 6

[69] Arash Vahdat and Jan Kautz. NVAE: A Deep Hierarchical Variational Autoencoder. In *Advances in Neural Information Processing Systems*, volume 33, 2020. 3

[70] Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based Generative Modeling in Latent Space. *arXiv preprint arXiv:2106.05931*, 2021. 1, 3

[71] Aaron van den Oord, Nal Kalchbrenner, Lasse Espeholt, koray kavukcuoglu, Oriol Vinyals, and Alex Graves. Conditional Image Generation with PixelCNN Decoders. *NeurIPS 29*, 2016. 2

[72] Aäron Van Den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel Recurrent Neural Networks. In *ICML*, 2016. 2

[73] Aaron van den Oord, Oriol Vinyals, and koray kavukcuoglu. Neural Discrete Representation Learning. *NeurIPS 30*, 2017. 1, 3

[74] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In *Advances in Neural Information Processing Systems*, volume 30, 2017. 2, 3, 4, 5

[75] Alex Wang and Kyunghyun Cho. BERT has a Mouth, and it Must Speak: BERT as a Markov Random Field Language Model. In *NeuralGen*, 2019. 2

[76] Daniel Watson, Jonathan Ho, Mohammad Norouzi, and William Chan. Learning to Efficiently Sample from Diffusion Probabilistic Models. *arXiv preprint arXiv:2106.03802*, 2021. 7

[77] Zhisheng Xiao, Karsten Kreis, Jan Kautz, and Arash Vahdat. VAEBM: A Symbiosis between Variational Autoencoders and Energy-based Models. In *International Conference on Learning Representations*, 2021. 1, 3

[78] 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. *arXiv preprint arXiv:1506.03365*, 2015. 5

[79] Ning Yu, Connelly Barnes, Eli Shechtman, Sohrab Amirghodsi, and Michal Lukac. Texture Mixer: A Network for Controllable Synthesis and Interpolation of Texture. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12164–12173, 2019. 1

[80] 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*, pages 586–595, 2018. 4, 13

[81] Shengyu Zhao, Zhijian Liu, Ji Lin, Jun-Yan Zhu, and Song Han. Differentiable Augmentation for Data-Efficient GAN Training. In *Advances in Neural Information Processing Systems*, volume 33, 2020. 4, 7

# Supplementary Material

The supplementary material for this work is divided into the following sections: Section A describes the architectures and hyperparameters for the experiments presented in the main paper; Section B illustrates the connection between our proposed ELBO reweighting and the true ELBO; Section C gives nearest neighbour examples to demonstrate generalisation; and finally, Section D contains a large number of uncurated samples.

## A. Implementation Details

We perform all experiments on a single NVIDIA RTX 2080 Ti with 11GB of VRAM using automatic mixed precision when possible. As mentioned in the main paper, we use the same VQGAN architecture as used by Esser et al. [18] which for $256 \times 256$ images downsamples to features of size $16 \times 16 \times 256$, and quantizes using a codebook with 1024 entries. Attention layers are applied within both the encoder and decoder on the lowest resolutions to aggregate context across the entire image. Models are optimised using the Adam optimiser [39] using a batch size of 4 and learning rate of $1.8 \times 10^{-5}$. For the differentiable augmentations we randomly change the brightness, saturation, and contrast, as well as randomly translate images. The datasets we use are both publically accessible, with FFHQ available under the Creative Commons BY 4.0 licence. LSUN models are trained for 2.2M steps while the FFHQ model is trained for 1.4M steps.

For the absorbing diffusion model we use a scaled down 80M parameter version of GPT-2 [57] consisting of 24 layers, where each attention layer has 8 heads, each 64D. The same architecture is used for experiments with the autoregressive model. Autoregressive models' training are stopped based on the best validation loss. We also stop training the absorbing diffusion models based on validation ELBO, however, on the LSUN datasets we found that it always improved or remained consistent throughout training so each model was trained for 2M steps.

**Codebook Collapse** One issue with vector quantized methods is codebook collapse, where some codes fall out of use which limits the potential expressivity of the model. We found this to occur across all datasets with often a fraction of the codes in use. We experimented with different quantization schemes such as gumbel softmax, different initalisation schemes such as k-means, and 'code recycling', where codes out of use are reset to an in use code. In all of these cases, we found the reconstruction quality to be comparable or worse so stuck with the argmax quantisation scheme used by Esser et al. [18].

**Precision, Recall, Density, and Coverage** To compute these measures we use the official code releases and pretrained weights in all cases except Taming Transformers on the LSUN datasets where weights were not available; in this case we reproduced results as close as possible with the hardware available, training the VQGANs and autoregressive models with the same hyperparameters used for the rest of our experiments. Following Nash et al. [51] we use the standard 2048D InceptionV3 features, which are also used to compute FID. The measures are computed using the code provided by Naeem et al. [50].

## B. Reweighted ELBO

In Section 3.2 we propose re-weighting the ELBO of the absorbing diffusion model so that the individual loss at each time step is multiplied by $\frac{T-t+1}{T}$ rather than $1/t$. In this section we justify the correctness of this re-weighting by showing it is equivalent to minimising the difference to a forward process that does not have access to $\boldsymbol{x}_t$. As such, the loss takes into account the difficulty of denoising steps and re-weights them down accordingly. In this case, the loss at time step $t$ can be written as

$$
\begin{aligned}
\mathcal{L}_t &= D_{\mathrm{KL}}(q(\boldsymbol{x}_{t-1}|\boldsymbol{x}_0)||p(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t)) \\
&= \sum_i \sum_j q([\boldsymbol{x}_{t-1}]_{i,j}|\boldsymbol{x}_0) \log \frac{q([\boldsymbol{x}_{t-1}]_{i,j}|\boldsymbol{x}_0)}{p([\boldsymbol{x}_{t-1}]_{i,j}|\boldsymbol{x}_t)},
\end{aligned} \quad (10)
$$

where the first summation sums over latent coordinates $i$, and the second summation sums over the probabilities of each code $j$. For the absorbing diffusion case where tokens in $\boldsymbol{x}_t$ are masked independently and uniformly with probability $\frac{t}{T}$, this posterior is defined as

$$
\begin{aligned}
&q([\boldsymbol{x}_{t-1}]_i = a|\boldsymbol{x}_0) \\
&= \begin{cases}
1 - \frac{t-1}{T}, & \text{if } a = [\boldsymbol{x}_0]_i \text{ and } [\boldsymbol{x}_t]_i = m. \\
\frac{t-1}{T}, & \text{if } a = m \text{ and } [\boldsymbol{x}_t]_i = m. \\
1, & \text{if } a = [\boldsymbol{x}_0]_i \text{ and } [\boldsymbol{x}_t]_i = [\boldsymbol{x}_0]_i. \\
0, & \text{otherwise.}
\end{cases}
\end{aligned} \quad (11)
$$

The reverse process remains defined in the same way as the standard reverse process:

$$
\begin{aligned}
&p([\boldsymbol{x}_{t-1}]_i = a|\boldsymbol{x}_t) \\
&= \begin{cases}
\frac{1}{t} p_\theta([\boldsymbol{x}_0]_i|\boldsymbol{x}_t), & \text{if } a = [\boldsymbol{x}_0]_i \text{ and } [\boldsymbol{x}_t]_i = m. \\
1 - \frac{1}{t}, & \text{if } a = m \text{ and } [\boldsymbol{x}_t]_i = m. \\
1, & \text{if } a = [\boldsymbol{x}_0]_i \text{ and } [\boldsymbol{x}_t]_i = [\boldsymbol{x}_0]_i.
\end{cases}
\end{aligned} \quad (12)
$$

Substituting these definitions into Equation (10), the loss can be simplified to Equation (13); by extracting the constants into a single term out of the sum, $C$, the loss can be

further simplified to obtain Equation (14), which is equivalent to our proposed reweighted ELBO Equation (10),

$$\mathcal{L}_t = \sum_i \left[ 1 \log \frac{1}{1} + \frac{t-1}{T} \log \frac{\frac{t-1}{T}}{1 - \frac{1}{t}} \right.$$
$$\left. + \left( 1 - \frac{t-1}{T} \log \frac{1 - \frac{t-1}{T}}{\frac{1}{t} p_\theta([\boldsymbol{x}_0]_i | \boldsymbol{x}_t)} \right) \right], \quad (13)$$

$$= C - \sum_i \left[ \frac{T - t - 1}{T} \log p_\theta([\boldsymbol{x}_0]_i | \boldsymbol{x}_t) \right]. \quad (14)$$

## C. Nearest Neighbours

When training generative models, being able to detect overfitting is key to ensure the data distribution is well modelled. Overfitting is not detected by popular metrics such as FID, making overfitting difficult to identify in approaches such as GANs. With our approach we are able to approximate the ELBO on a validation set making it simple to prevent overfitting. In this section we demonstrate that our approach is not overfit by providing nearest neighbour images from the training dataset to samples from our model, measured using LPIPS [80].

## D. Additional Samples

Figure 11 contains unconditional samples with resolutions larger than observed in the training data from a model trained on LSUN Bedroom. In Figures 12 to 14 additional samples from our models are visualised on LSUN Churches, LSUN Bedroom, and FFHQ.

Figure 8. Nearest neighbours for a model trained on LSUN Churches based on LPIPS distance. The left column contains samples from our model and the right column contains the nearest neighbours in the training set (increasing in distance from left to right).



Figure 9. Nearest neighbours for a model trained on FFHQ based on LPIPS distance. The left column contains samples from our model and the right column contains the nearest neighbours in the training set (increasing in distance from left to right).

Figure 10. Nearest neighbours for a model trained on LSUN Bedroom based on LPIPS distance. The left column contains samples from our model and the right column contains the nearest neighbours in the training set (increasing in distance from left to right).

Figure 11. Unconditional samples from a model trained on LSUN Bedroom larger than images in the training dataset.

Figure 12. Non-cherry picked, $t = 0.9$, 256x256 LSUN Churches samples.

Figure 13. Non-cherry picked, $t = 0.85$, 256x256 FFHQ samples.

Figure 14. Non-cherry picked, $t = 0.9$, 256x256 LSUN Bedroom samples.