
Anytime Sampling for Autoregressive Models via Ordered Autoencoding

Anonymous Author(s)

Affiliation

Address

email

Abstract

Autoregressive models are widely used for tasks such as image and audio generation. Sampling from autoregressive models, however, can be slow, involving a computation that is sequential in nature and typically scales linearly with respect to the data dimension. To address this difficulty, we propose a new family of autoregressive models that enables anytime sampling. Inspired by Principal Component Analysis (PCA), we learn a structured representation space where dimensions are ordered based on their importance with respect to reconstruction. Using an autoregressive model in this latent space, we trade off sample quality for computational efficiency by truncating the generation process before decoding into the original data space. Experimentally, we demonstrate in several image and audio generation tasks that sample quality degrades gracefully as we reduce the computational budget for sampling. The approach suffers almost no loss in sample quality (measured by FID) using only 60% to 80% of all latent dimensions for image data.

1 Introduction

Autoregressive models are a prominent approach to data generation, and have been widely used to produce high quality samples of images (Oord et al., 2016b; Salimans et al., 2017; Menick & Kalchbrenner, 2018), audio (Oord et al., 2016a), video (Kalchbrenner et al., 2017) and text (Kalchbrenner et al., 2016; Radford et al., 2019). These models represent a joint distribution as a product of (simpler) conditionals, and sampling requires iterating over all these conditional distributions in a certain order. As a result, the sampling process of autoregressive models can be slow. Due to the sequential nature of this process, the computational cost will grow at least linearly with respect to the number of conditional distributions, which is typically equal to the data dimension.

Although caching techniques have been developed to speed up generation (Ramachandran et al., 2017; Guo et al., 2017), the high cost of sampling limits their applicability in many scenarios. For example, when running on multiple devices with different computational resources, we may wish to trade off sample quality for faster generation based on the computing power available on each device. Currently, a separate model must be trained for each device (*i.e.*, computational budget) in order to trade off sample quality for faster generation, and there is no way to control this trade-off on the fly to accommodate instantaneous resource availability at time-of-deployment.

To address this difficulty, we consider the novel task of *adaptive* autoregressive generation under *computational constraints*. We seek to build a *single model* that can automatically trade-off sample quality versus computational cost via *anytime sampling*, *i.e.*, where the sampling process may be terminated anytime (*e.g.*, because of exhausted computational budget) to yield a *complete* sample whose sample quality decays with the earliness of termination.

36 In particular, we take advantage of a generalization of Principal Components Analysis (PCA) pro-
 37 posed by Rippel et al. (2014), which learns an ordered representations induced by a structured
 38 application of dropout to the representations learned by an autoencoder. Such a representation
 39 encodes raw data into a latent space where dimensions are sorted based on their importance for
 40 reconstruction. Autoregressive modeling is then applied in the ordered representation space instead.
 41 This approach enables a natural trade-off between quality and computation by truncating the length
 42 of the representations: When running on devices with high computational capacity, we can afford
 43 to generate the full representation and decode it to obtain a high quality sample; when on a tighter
 44 computational budget, we can generate only the first few dimensions of the representation and decode
 45 it to a sample whose quality degrades smoothly with truncation. Because decoding is usually fast and
 46 the main computation bottleneck lies on the autoregressive part, the run-time grows proportionally
 47 relative to the number of sampled latent dimensions.
 48 Through experiments, we show that our autoregressive models are capable of trading off sample
 49 quality and inference speed. When training autoregressive models on the latent space given by our
 50 encoder, we witness little degradation of image sample quality using only around 60% to 80% of all
 51 latent codes, as measured by Fréchet Inception Distance (Heusel et al., 2017b) on CIFAR-10 and
 52 CelebA. Compared to standard autoregressive models, our approach allows the sample quality to
 53 degrade gracefully as we reduce the computational budget for sampling. We also observe that on
 54 the VCTK audio dataset (Veaux et al., 2017), our autoregressive model is able to generate the low
 55 frequency features first, then gradually refine the waveforms with higher frequency components as
 56 we increase the number of sampled latent dimensions.

57 2 Background

58 **Autoregressive Models** Autoregressive models define a probability distribution over data points
 59 $\mathbf{x} \in \mathbb{R}^D$ by factorizing the joint probability distribution as a product of univariate conditional
 60 distributions with the chain rule. Using p_θ to denote the distribution of the model, we have:

$$p_\theta(\mathbf{x}) = \prod_{i=1}^D p_\theta(x_i | x_1, \dots, x_{i-1}) \quad (1)$$

61 The model is trained by maximizing the likelihood:

$$\mathcal{L} = \mathbb{E}_{p_d(\mathbf{x})}[\log p_\theta(\mathbf{x})], \quad (2)$$

62 where $p_d(\mathbf{x})$ represents the data distribution.

63 Different autoregressive models adopt different orderings of input dimensions and parameterize
 64 the conditional probability $p_\theta(x_i | x_1, \dots, x_{i-1}), i = 1, \dots, D$ in different ways. Most architec-
 65 tures over images order the variables x_1, \dots, x_D of image \mathbf{x} in *raster scan order* (i.e., left-to-right
 66 then top-to-bottom). Popular autoregressive architectures include MADE (Papamakarios et al.,
 67 2017), PixelCNN (Oord et al., 2016b; van den Oord et al., 2016; Salimans et al., 2017) and Trans-
 68 former (Vaswani et al., 2017), where they respectively use masked linear layers, convolutional layers
 69 and self-attention blocks to ensure that the output corresponding to $p_\theta(x_i | x_1, \dots, x_{i-1})$ is oblivious
 70 of x_i, x_{i+1}, \dots, x_D .

71 **Cost of Sampling** During training, we can evaluate autoregressive models efficiently because
 72 x_1, \dots, x_D are provided by data and all conditionals $p(x_i | x_1, \dots, x_{i-1})$ can be computed in
 73 parallel. In contrast, sampling from autoregressive models is an inherently sequential process and
 74 cannot be easily accelerated by parallel computing: we first need to sample x_1 , after which we sample
 75 x_2 from $p_\theta(x_2 | x_1)$ and so on—the i -th variable x_i can only be obtained after we have already
 76 computed x_1, \dots, x_{i-1} . Thus, the run-time of autoregressive generation grows *at least linearly* with
 77 respect to the length of a sample. In practice, the sample length D can be more than hundreds of
 78 thousands for real-world image and audio data. This poses a major challenge to fast autoregressive
 79 generation on a small computing budget.

80 3 Anytime Sampling with Ordered Autoencoders

81 Our goal is to circumvent this linear time complexity of autoregressive models by pushing the task of
 82 autoregressive modeling from the original data space (e.g., pixel space) into an ordered representation

space. In doing so, we develop a new class of autoregressive models where premature truncation of the autoregressive sampling process leads to the generation of a *lower quality* sample instead of an *incomplete* sample. In this section, we shall first describe the learning of the ordered representation space via the use of an *ordered autoencoder*. We then describe how to achieve anytime sampling with ordered autoencoders.

3.1 Ordered Autoencoders

Consider an autoencoder that encodes an input $\mathbf{x} \in \mathbb{R}^D$ to a code $\mathbf{z} \in \mathbb{R}^K$. Let $\mathbf{z} = e_\theta(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}^K$ be the encoder parameterized by θ and $\mathbf{x}' = d_\phi(\mathbf{z}) : \mathbb{R}^K \rightarrow \mathbb{R}^D$ be the decoder parameterized by ϕ . We define $e_\theta(\cdot)_{\leq i} : \mathbf{x} \in \mathbb{R}^D \mapsto (z_1, z_2, \dots, z_i, 0, \dots, 0)^T \in \mathbb{R}^K$, which truncates the representation to the first i dimensions of the encoding $\mathbf{z} = e_\theta(\mathbf{x})$, masking out the remainder of the dimensions with a zero value. We define the ordered autoencoder objective as

$$\frac{1}{N} \sum_{i=1}^N \frac{1}{K} \sum_{j=1}^K \|\mathbf{x}_i - d_\phi(e_\theta(\mathbf{x}_i)_{\leq j})\|_2^2. \quad (3)$$

We note that Eq. (3) is equivalent to Rippel et al. (2014)'s nested dropout formulation using a uniform sampling of possible truncations. Moreover, when the encoder/decoder pair is constrained to be a pair of orthogonal matrices up to a transpose, then the optimal solution in Eq. (3) recovers PCA.

3.1.1 Theoretical Analysis

Rippel et al. (2014)'s analysis of the ordered autoencoder is limited to linear/sigmoid encoder and a linear decoder. In this section, we extend the analysis to general autoencoder architectures by employing an information-theoretic framework to analyze the importance of the i -th latent code to reconstruction for ordered autoencoders. We first reframe our problem from a probabilistic perspective. In lieu of using deterministic autoencoders, we assume that both the encoder and decoder are stochastic functions. In particular, we let $q_{e_\theta}(\mathbf{z} | \mathbf{x})$ be a probability distribution over $\mathbf{z} \in \mathbb{R}^K$ conditioned on input \mathbf{x} , and similarly let $p_{d_\phi}(\mathbf{x} | \mathbf{z})$ be the stochastic counterpart to $d_\phi(\mathbf{z})$. We then use $q_{e_\theta}(\mathbf{z} | \mathbf{x})_{\leq i}$ to denote the distribution of $(z_1, z_2, \dots, z_i, 0, \dots, 0)^T \in \mathbb{R}^K$, where $\mathbf{z} \sim q_{e_\theta}(\mathbf{z} | \mathbf{x})$, and let $p_{d_\phi}(\mathbf{x} | \mathbf{z})_{\leq i}$ represent the distribution of $p_{d_\phi}(\mathbf{x} | (z_1, z_2, \dots, z_i, 0, \dots, 0)^T \in \mathbb{R}^K)$. We can modify Eq. (3) to have the following form:

$$\mathbb{E}_{x \sim p_d(\mathbf{x}), i \sim \mathcal{U}\{1, K\}} \mathbb{E}_{\mathbf{z} \sim q_{e_\theta}(\mathbf{z} | \mathbf{x})_{\leq i}} [-\log p_{d_\phi}(\mathbf{x} | \mathbf{z})_{\leq i}], \quad (4)$$

where $\mathcal{U}\{1, K\}$ denotes a uniform distribution over $\{1, 2, \dots, K\}$, and $p_d(\mathbf{x})$ represents the data distribution. We can choose both the encoder and decoder to be fully factorized Gaussian distributions with a fixed variance σ^2 , then Eq. (4) can be simplified to

$$\mathbb{E}_{p_d(\mathbf{x})} \left[\frac{1}{K} \sum_{i=1}^K \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(e_\theta(\mathbf{x})_{\leq i}; \sigma^2)} \left[\frac{1}{2\sigma^2} \|\mathbf{x} - d_\phi(\mathbf{z})_{\leq i}\|_2^2 \right] \right].$$

The stochastic encoder and decoder in this case will become deterministic when $\sigma \rightarrow 0$, and the above equation will yield the same encoder/decoder pair as Eq. (3) when $\sigma \rightarrow 0$ and $N \rightarrow \infty$.

The optimal encoders and decoders that minimize Eq. (4) satisfy the following property.

Theorem 1. *Let \mathbf{x} denote the input random variable. Assuming both the encoder and decoder are optimal in terms of minimizing Eq. (4), and $\forall i \in 2, \dots, K, z_{i-1} \perp z_i | \mathbf{x}$, we have*

$$\forall i \in \{2, \dots, K\} : I(z_i; \mathbf{x} | z_{\leq i-1}) \leq I(z_{i-1}; \mathbf{x} | z_{\leq i-2}),$$

where $z_{\leq i}$ denotes (z_1, z_2, \dots, z_i) .

We defer the proof to Appendix A.1. The assumption $z_{i-1} \perp z_i | \mathbf{x}$ holds whenever the encoder $q_\theta(\mathbf{z} | \mathbf{x})$ is a factorized distribution, which is a common choice in variational autoencoders (Kingma & Welling, 2013), and we use $I(\mathbf{a}; \mathbf{b} | \mathbf{c})$ to denote the mutual information between random variables \mathbf{a} and \mathbf{b} conditioned on \mathbf{c} . Intuitively, the above theorem states that for optimal encoders and decoders that minimize Eq. (4), one can extract less additional information about the raw input as the code gets longer. Therefore, there exists a natural ordering among different dimensions of the code based on the additional information they can provide for reconstructing the inputs.

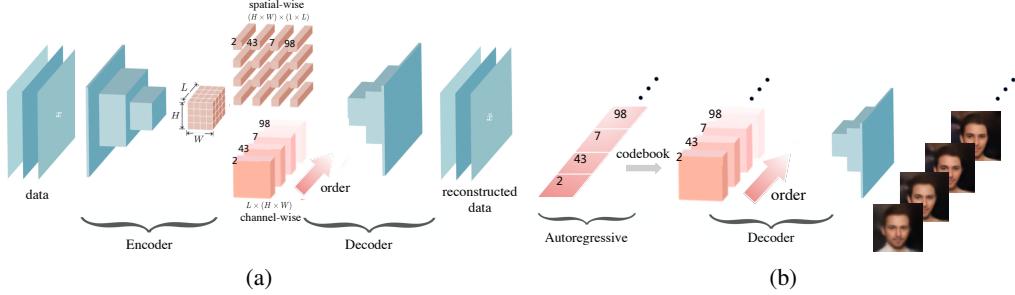


Figure 1: (a) Spatial-wise quantization vs. channel-wise quantization. (b) Anytime sampling for OVQ-VAE.

124 3.2 Anytime Sampling

125 Once we have learned an ordered autoencoder, we then train an autoregressive model on the full
 126 length codes in the ordered representation space, to adopt the ex-post density estimation [Ghosh et al.
 127 (2020)]. For each input \mathbf{x}_i in a dataset $\mathbf{x}_1, \dots, \mathbf{x}_N$, we feed it to the encoder to get $\mathbf{z}_i = e_\theta(\mathbf{x}_i)$. The
 128 resulting codes $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N$ are used as training data. After training both the ordered autoencoder
 129 and autoregressive model, we can perform anytime sampling on a large spectrum of computing
 130 budgets. Suppose for example we can afford to generate T code dimensions from the autoregressive
 131 model, denoted as $\mathbf{z}_{\leq T} \in \mathbb{R}^T$. We can simply zero-pad it to get $(z_1, z_2, \dots, z_T, 0, \dots, 0)^T \in \mathbb{R}^K$
 132 and decode it to get a complete sample. Unlike the autoregressive part, the decoder has access to
 133 all dimensions of the latent code at the same time and can decode in parallel. On modern GPUs,
 134 the code length has minimal effect on the run-time of decoding, as long as the decoder is not itself
 135 autoregressive (see empirical verifications in Section 5.2.2).

136 4 Ordered VQ-VAE

137 In this section, we apply the ordered autoencoder framework to the vector quantized variational
 138 autoencoder (VQ-VAE) and its extension [van den Oord et al., 2017; Razavi et al., 2019]. Since
 139 these models are quantized autoencoders paired with a latent autoregressive model, they admit a
 140 natural extension to *ordered* VQ-VAEs (OVQ-VAEs) under our framework—a new family of VQ-
 141 VAE models capable of anytime sampling. Below, we begin by describing the VQ-VAE, and then
 142 highlight two key design choices (ordered discrete codes and channel-wise quantization) critical for
 143 OVQ-VAEs.

144 4.1 VQ-VAE

145 To construct a VQ-VAE with code length of K discrete latent variables, the encoder must first
 146 map the raw input \mathbf{x} to a continuous representation $\mathbf{z}_e = e_\theta(\mathbf{x}) \in \mathbb{R}^{K \times D}$, before feeding it to a
 147 vector-valued quantization function $q : \mathbb{R}^{K \times D} \rightarrow \{1, 2, \dots, C\}^K$ defined as

$$q(\mathbf{z}_e)_j = \arg \min_{i \in \{1, \dots, C\}} \|\mathbf{e}_i - \mathbf{z}_{e,j}\|_2,$$

148 where $q(\mathbf{z}_e)_j \in \{1, 2, \dots, C\}$ denotes the j -th component of the vector-valued function $q(\mathbf{z}_e)$,
 149 $\mathbf{z}_{e,j} \in \mathbb{R}^D$ denotes the j -th row of \mathbf{z}_e , and \mathbf{e}_i denotes the i -th row of the embedding matrix
 150 $\mathbf{E} \in \mathbb{R}^{C \times D}$. Next, we view $q(\mathbf{z}_e)$ as a sequence of indices and use them to look up embedding
 151 vectors from the codebook \mathbf{E} . This yields a latent representation $\mathbf{z}_d \in \mathbb{R}^{K \times D}$, given by $\mathbf{z}_{d,j} = \mathbf{e}_{q(\mathbf{z}_e)_j}$,
 152 where $\mathbf{z}_{d,j} \in \mathbb{R}^D$ denotes the j -th row of \mathbf{z}_d . Finally, we can decode \mathbf{z}_d to obtain the reconstruction
 153 $d_\phi(\mathbf{z}_d)$. This procedure can be viewed as a regular autoencoder with a non-differentiable nonlinear
 154 function that maps each latent vector $\mathbf{z}_{e,j}$ to 1-of- K embedding vectors \mathbf{e}_i .

155 During training, we use the straight-through gradient estimator [Bengio et al., 2013] to propagate
 156 gradients through the quantization function, *i.e.*, gradients are directly copied from the decoder input
 157 \mathbf{z}_d to the encoder output \mathbf{z}_e . The loss function for training on a single data point \mathbf{x} is given by

$$\|d_\phi(\mathbf{z}_d) - \mathbf{x}\|_2^2 + \|\text{sg}[e_\theta(\mathbf{x})] - \mathbf{z}_d\|_F^2 + \beta \|e_\theta(\mathbf{x}) - \text{sg}[\mathbf{z}_d]\|_F^2, \quad (5)$$

158 where sg stands for the `stop_gradient` operator, and β is a hyper-parameter ranging from 0.1 to
 159 2.0. The first term of Eq. (5) is the standard reconstruction loss, the second term is for embedding
 160 learning while the third term is for training stability (van den Oord et al., 2017). Samples from a
 161 VQ-VAE can be produced by first training an autoregressive model on its latent space, followed by
 162 decoding samples from the autoregressive model into the raw data space.

163 4.2 Ordered Discrete Codes

164 Since the VQ-VAE outputs a sequence of discrete latent codes $q(\mathbf{z}_e)$, we wish to impose an ordering
 165 that prioritizes the code dimensions based on importance to reconstruction. In analogy to Eq. (3),
 166 we can modify the reconstruction error term in Eq. (5) to learn ordered latent representations. The
 167 modified loss function is an order-inducing objective given by

$$\frac{1}{K} \sum_{i=1}^K [\|d_\phi(\mathbf{z}_{d,\leq i}) - \mathbf{x}\|_2^2 + \|\text{sg}[e_\theta(\mathbf{x})_{\leq i}] - \mathbf{z}_{d,\leq i}\|_F^2 + \beta \|e_\theta(\mathbf{x})_{\leq i} - \text{sg}[\mathbf{z}_{d,\leq i}]\|_F^2], \quad (6)$$

168 where $e_\theta(\mathbf{x})_{\leq i}$ and $\mathbf{z}_{d,\leq i}$ denote the results of keeping the top i rows of $e_\theta(\mathbf{x})$ and \mathbf{z}_d and then masking
 169 out the remainder rows with zero vectors. We uniformly sample the masking index $i \sim \mathcal{U}\{1, K\}$ to
 170 approximate the average in Eq. (6) when K is large. This results in the learning of ordered discrete
 171 latent variables, which can then be paired with a latent autoregressive model for anytime sampling.

172 4.3 Channel-Wise Quantization

173 In Section 4.1 we assume the encoder output to be a $K \times D$ matrix (*i.e.*, $\mathbf{z}_e \in \mathbb{R}^{K \times D}$). In practice,
 174 the output can have various sizes depending on the encoder network, and we need to reshape it
 175 to a two-dimensional matrix. For example, when encoding images, the encoder is typically a 2D
 176 convolutional neural network (CNN) whose output is a 3D latent feature map of size $L \times H \times W$.
 177 Here L , H , and W stand for the channel, height, and width of the feature maps. We discuss below
 178 how the reshaping procedure can significantly impact the performance of anytime sampling and
 179 propose a reshaping procedure that facilitates high-performance anytime sampling.

180 Consider convolutional encoders on image data, where the output feature map has a size of $L \times H \times W$.
 181 The most common way of reshaping

182 this 3D feature map, as in van den
 183 Oord et al. (2017), is to let $H \times W$ be
 184 the code length, and let the number of
 185 channels L be the size of embedding
 186 vectors, *i.e.*, $K = H \times W$ and $D = L$. We call this pattern *spatial-wise*
 187 *quantization*, as each spatial location
 188 in the feature map corresponds to one
 189 code dimension and will be quantized
 190 separately. Since the code dimensions
 191 correspond to spatial locations of the
 192 feature map, they encode local features due to a limited receptive field. This is detrimental to anytime
 193 sampling, because early dimensions cannot capture the global information needed for reconstructing
 194 the entire image. We demonstrate this in Fig. 2, which shows that VQ-VAE with spatial-wise
 195 quantization is only able to reconstruct the top rows of an image with $1/4$ of the code length.

196 To address this issue, we propose *channel-wise quantization*, where each channel of the feature map
 197 is viewed as one code dimension and quantized separately (see Fig. I(a) for visual comparison of
 198 spatial-wise and channel-wise quantization). Specifically, the code length is L (*i.e.*, $K = L$), and
 199 the size of the embedding vectors is $H \times W$ (*i.e.*, $D = H \times W$). In this case, one code dimension
 200 includes all spatial locations in the feature map and can capture global information better. As shown
 201 in the right panel of Fig. 2, channel-wise quantization clearly outperforms spatial-wise quantization
 202 for anytime sampling. We use channel-wise quantization in all subsequent experiments.

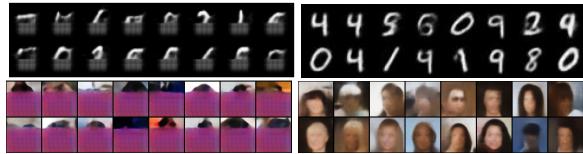


Figure 2: MNIST (**top**) and CelebA (**bottom**) samples generated with $1/4$ of the code length. **Left:** Spatial-wise quantization. **Right:** Channel-wise quantization.

204 5 Experiments

205 In our experiments, we focus on anytime sampling for autoregressive models trained on the latent
 206 space of OVQ-VAEs. We first verify that our learning objectives Eq. (3), Eq. (6) are effective at

207 inducing ordered representations. Next, we demonstrate on several image and audio datasets that our
 208 models achieve comparable sample quality to regular VQ-VAEs, while allowing graceful trade-off
 209 between sample quality and computation time via anytime sampling. Due to the page limit, we defer
 210 the result of audio generation to Appendix B and provide additional experimental details and results
 211 in Appendix C and E.

212 5.1 Ordered Versus Unordered Codes

213 Our proposed ordered autoencoder framework learns an ordered encoding that is in contrast to
 214 the encoding learned by a standard au-
 215 toencoder (which we shall refer to as
 216 unordered). In addition to the theo-
 217 retical analysis in Section 3.1.1, we
 218 provide further empirical analysis to
 219 characterize the difference between or-
 220 dered and unordered codes. In particu-
 221 lar, we compare the importance of the
 222 i -th code—as measured by the reduc-
 223 tion in reconstruction error $\Delta(i)$ —for
 224 PCA, standard (unordered) VQ-VAE,
 225 and ordered VQ-VAE. For VQ-VAE
 226 and ordered VQ-VAE, we define the
 227 reduction in reconstruction error $\Delta(i)$
 228 for a data point \mathbf{x} as

$$\Delta_{\mathbf{x}}(i) \triangleq \|d_{\phi}(\mathbf{z}_{d,\leq i-1}) - \mathbf{x}\|_F^2 - \|d_{\phi}(\mathbf{z}_{d,\leq i}) - \mathbf{x}\|_F^2 \quad (7)$$

229 Averaging $\Delta_{\mathbf{x}}(i)$ over the entire dataset thus yields $\Delta(i)$. Similarly we define $\Delta(i)$ as the reduction
 230 on reconstruction error of the entire dataset, when adding the i -th principal component for PCA.

231 Fig. 3(a) shows the $\Delta(i)$'s of the three models on the CIFAR-10. Since PCA and ordered VQ-VAE
 232 both learn an ordered encoding, their $\Delta(i)$'s decay gradually as i increases. In contrast, the standard
 233 VQ-VAE with an unordered encoding exhibits a highly irregular $\Delta(i)$, indicating no meaningful
 234 ordering of the dimensions.

235 Fig. 3(b) further shows how the reconstruction error decreases as a function of the truncated code
 236 length for the three models. Although unordered VQ-VAE and ordered VQ-VAE achieve similar
 237 reconstruction errors for sufficiently large code lengths, it is evident that an ordered encoding achieves
 238 significantly better reconstructions when the code length is aggressively truncated. When sufficiently
 239 truncated, we observe even PCA outperforms unordered VQ-VAE despite the latter being a more
 240 expressive model. In contrast, ordered VQ-VAE achieves superior reconstructions compared to PCA
 241 and unordered VQ-VAE across all truncation lengths.

242 5.2 Image Generation

243 We test the performance of anytime sampling using OVQ-VAEs on several image datasets. We
 244 compare our approach to two baselines. One is the original VQ-VAE model proposed by van den
 245 Oord et al. (2017) without anytime sampling. The other is using anytime sampling with unorderd
 246 VQ-VAEs, where the models have the same architectures as ours but are trained by minimizing
 247 Eq. (5). We empirically verify that 1) we are able to generate high quality image samples; 2) image
 248 quality degrades gracefully as we reduce the sampled code length for anytime sampling; and 3)
 249 anytime sampling improves the inference speed compared to naïve sampling of original VQ-VAEs.

250 We evaluate the model performance on the MNIST, CIFAR-10 (Krizhevsky, 2009) and CelebA (Liu
 251 et al., 2014) datasets. For CelebA, the images are resized to 64×64 . All pixel values are scaled to
 252 the range $[0, 1]$. We borrow the model architectures and optimizers from van den Oord et al. (2017)
 253 and adopt some minor modifications. We empirically found that the residual blocks in the original
 254 VQ-VAE architecture can lead to unstable training when applying channel-wise quantization, so we
 255 replace these residual blocks with convolutional layers, while keeping the number of parameters
 256 approximately fixed. The full code length and the codebook size are 16 and 126 for MNIST, 70 and
 257 1000 for CIFAR-10, and 100 and 500 for CelebA respectively. We train a Transformer (Vaswani
 258 et al., 2017) on our VQ-VAEs, as opposed to the PixelCNN model used in van den Oord et al. (2017).

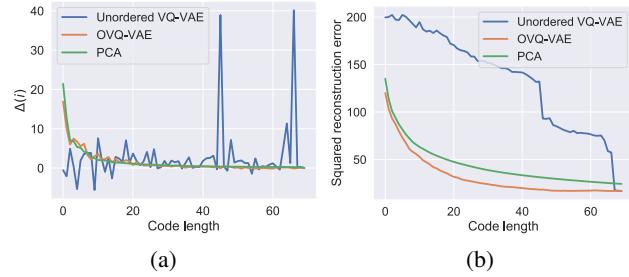


Figure 3: Ordered vs. unordered codes on CIFAR-10.

259 Transformers are arguably more suitable for channel-wise quantization, since there are no spatial
 260 relations among different code dimensions that can be leveraged by convolutional models (such as
 261 PixelCNNs).

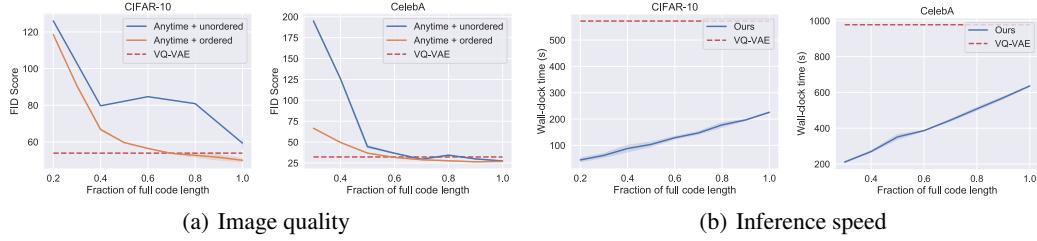


Figure 4: (a) FID scores for anytime sampling using various code lengths. (b) Inference speed of anytime sampling with different code lengths.

262 5.2.1 Image Quality

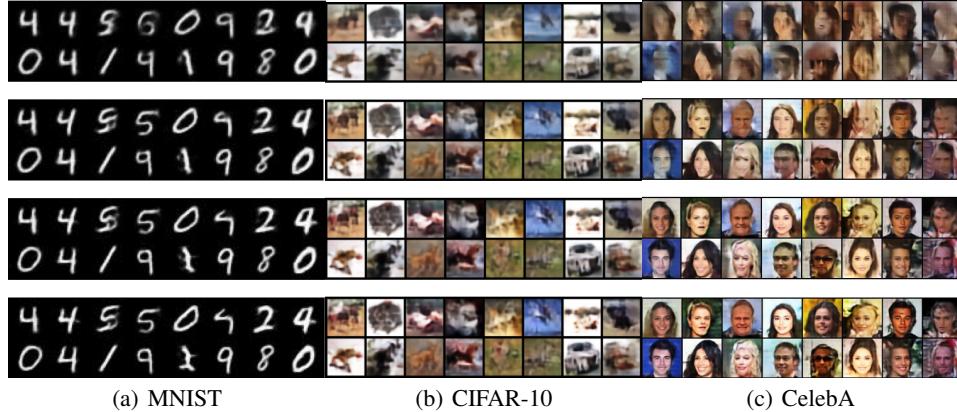


Figure 5: Anytime sampling with 0.25/0.5/0.75/1.0 (**top to bottom**) fractions of full code length on MNIST (**left**), CIFAR-10 (**middle**) and CelebA (**right**) datasets.

263 In Fig. 4(a), we report FID (Heusel et al., 2017b) scores (a popular metric on image quality) on
 264 CIFAR-10 and CelebA when performing anytime sampling for ordered versus unordered VQ-VAE.
 265 As a reference, we also report the FID scores when using the original VQ-VAE model (with residual
 266 blocks and spatial-wise quantization) sampled at the full code length (van den Oord et al., 2017). Our
 267 main finding is that OVQ-VAE achieves a better FID score than unordered VQ-VAE at all fractional
 268 code lengths (ranging from 20% to 100%); in other words, OVQ-VAE achieves strictly superior
 269 anytime sampling performance compared to unordered VQ-VAE on both CIFAR-10 and CelebA.
 270 In Fig. 5 (more in Appendix E.1), we visualize the sample quality degradation as a function of
 271 fractional code length when sampling from the OVQ-VAE. We observe a consistent increase in
 272 sample quality as we increased the fractional code length. In particular, we observe the model to
 273 initially generate a global structure of an image and then gradually fill in local details. We further show
 274 in Appendix D that, samples sharing the highest priority latent code have similar global structure.

275 5.2.2 Inference Speed

276 We compare the inference speed of our approach vs. the original VQ-VAE model by the wall-clock
 277 time needed for sampling. We respectively measure the time of generating 50000 and 100000 images
 278 on CIFAR-10 and CelebA datasets, with a batch size of 100. All samples are produced on a single
 279 NVIDIA TITAN Xp GPU.

280 Fig. 4(b) shows that the time needed for anytime sampling increases almost linearly with respect to
 281 the sampled code length. This supports our argument in Section 3.2 that the run-time of decoding is

282 negligible compared to the autoregressive component. Indeed, the decoder took around 23 seconds to
283 generate all samples for CelebA, whereas the sampling time of the autoregressive model was around
284 610 seconds—over an order of magnitude larger. Moreover, since we can achieve roughly the highest
285 sample quality with only 60% of the full code length on CelebA, anytime sampling can save around
286 40% run-time compared to naïve sampling without hurting sample quality.

287 In addition, our method is faster than the original VQ-VAE even when sampling the full code length,
288 without compromising sample quality (*cf.*, Section 5.2.1). This is because the Transformer model
289 we used is sufficiently shallower than the Gated PixelCNN (van den Oord et al., 2016) model in the
290 original VQ-VAE paper. Compared to PixelCNN++ (Salimans et al., 2017), an autoregressive model
291 on the raw pixel space, the sampling speed of our method can be an order of magnitude faster since
292 our autoregressive models are trained on the latent space with much lower dimensionality.

293 6 Related Work

294 Prior work has tackled the issue of slow autoregressive generation by improving implementations of
295 the generation algorithm. For example, the sampling speed of convolutional autoregressive models
296 can be improved substantially by caching hidden state computation (Ramachandran et al., 2017).
297 While such approaches provide substantial speedups in generation time, they are still at best linear
298 in the dimension of the sample space. van den Oord et al. (2018) improves the inference speed by
299 allowing parallel computing. Compared to our approach, they do not have the test-time adaptivity to
300 computational constraints. In contrast, we design methods that allow trade-offs between generation
301 speed and sample quality on the fly based on computational constraints, without model re-training.
302 For example, running apps need to accommodate to the real-time computational resources without
303 model re-training.

304 In order to enable anytime sampling, our method requires learning an ordered latent representation
305 of data by training ordered autoencoders. A similar method for training ordered encoders was
306 proposed in Rippel et al. (2014). Instead of uniform distribution over discrete codes in our method,
307 they adopted a geometric distribution over continuous codes during the training of the ordered
308 autoencoders. Because of the difference they require additional tricks such as unit sweeping and
309 adaptive regularization coefficients to stabilize the training, while our training is stable and more
310 suitable for large scale training. In addition, they only focus on fast retrieval and image compression.
311 By contrast, we further extend our approach to autoencoders with discrete latent codes (*e.g.*, VQ-
312 VAEs) and explore their applications in anytime sampling for autoregressive models. Another work
313 related to our approach is hierarchical nonlinear PCA (Scholz & Vigário, 2002). We generalize
314 their approach to latent spaces of arbitrary dimensionality, and leverage Monte Carlo estimations to
315 improve the efficiency when learning very high dimensional latent representations.

316 A line of works draw connections between PCA and the linear autoencoders. Kunin et al. (2019)
317 prove that the principal directions can be deduced from the critical points of L_2 regularized linear
318 autoencoders. Grover & Ermon (2019) show that the optimal measurement matrix recovers the
319 principal components and Pfau et al. (2019) use deep neural network to approximate eigenfunctions
320 of linear operators.

321 7 Conclusion

322 Sampling from autoregressive models is an expensive sequential process that can be intractable
323 when on a tight computing budget. To address this difficulty, we consider the novel task of adaptive
324 autoregressive sampling that can naturally trade-off computation with sample quality. Inspired by
325 PCA, we adopt ordered autoencoders, whose latent codes are prioritized based on their importance to
326 reconstruction. We show that it is possible to do anytime sampling for autoregressive models trained
327 on these ordered latent codes—we may stop the sequential sampling process at any step and still
328 obtain a complete sample of reasonable quality by decoding the partial codes.

329 With both theoretical arguments and empirical evidence, we show that ordered autoencoders can
330 induce a valid ordering that facilitates anytime sampling. Experimentally, we test our approach on
331 several image and audio datasets by pairing an ordered VQ-VAE (a powerful autoencoder architecture)
332 and a Transformer (an expressive autoregressive model) on the latent space. We demonstrate that
333 our samples suffer almost no loss of quality (as measured by FID scores) for images when using
334 only 60% to 80% of all code dimensions, and the sample quality degrades gracefully as we gradually
335 reduce the code length.

336 **Broader Impact Statement**

337 Generative models is a burgeoning research field with widespread implications for science and society.
338 Our work addresses the computational bottleneck in using autoregressive generative models for data
339 generation. Given the powerful expressive capabilities of autoegressive models, this aids many
340 downstream applications and in particular those that are limited by computational constraints and
341 require *fast, real-time data generation*.

342 Depending on the specific downstream usecase, the use of these models can have both positive and
343 negative outcomes. For example, autoregressive models with anytime sampling can be incorporated in
344 producing image/audio samples on-the-fly. However deepfake technology can take the advantage of it
345 and lead to social engineering scams. Autoregressive models with anytime sampling can also be used
346 for planning and reinforcement learning algorithms to enable real-time decision making in complex
347 environments. This further provides avenues for deployment in autonomous embodied artificial
348 intelligence systems, such as self-driving cars. On one hand, such systems can save lives in certain
349 scenarios e.g., error due to humans paying less attention to other traffic while driving can be reduced.
350 At the same time, the current systems are too brittle and susceptible to adversarial examples leading
351 to unanticipated failure. Besides the aforementioned limitations of such systems, there are also ethical
352 concerns in deploying autonomous systems, e.g., the trolley problem. Addressing the above concerns
353 requires further research in providing robustness guarantees for such model-based decision making
354 algorithms as well as close collaborations with researchers in socio-technical disciplines.

355 **References**

- 356 Bengio, Y., Léonard, N., and Courville, A. Estimating or propagating gradients through stochastic neurons for
357 conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- 358 Donahue, C., McAuley, J. J., and Puckette, M. Adversarial audio synthesis. In *ICLR*, 2018.
- 359 Ghosh, P., Sajjadi, M. S. M., Vergari, A., Black, M. J., and Schölkopf, B. From variational to deterministic
360 autoencoders. *ArXiv*, abs/1903.12436, 2020.
- 361 Grover, A. and Ermon, S. Uncertainty autoencoders: Learning compressed representations via variational
362 information maximization. *ArXiv*, abs/1812.10539, 2019.
- 363 Guo, P., Ni, X., Chen, X., and Ji, X. Fast pixelcnn: Based on network acceleration cache and partial generation
364 network. *2017 International Symposium on Intelligent Signal Processing and Communication Systems
(ISPACS)*, pp. 71–76, 2017.
- 366 Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. Gans trained by a two time-scale
367 update rule converge to a local nash equilibrium. In *NIPS*, 2017a.
- 368 Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. Gans trained by a two time-scale
369 update rule converge to a local nash equilibrium. In *Advances in neural information processing systems*, pp.
370 6626–6637, 2017b.
- 371 Kalchbrenner, N., Espeholt, L., Simonyan, K., Oord, A. v. d., Graves, A., and Kavukcuoglu, K. Neural machine
372 translation in linear time. *arXiv preprint arXiv:1610.10099*, 2016.
- 373 Kalchbrenner, N., van den Oord, A., Simonyan, K., Danihelka, I., Vinyals, O., Graves, A., and Kavukcuoglu, K.
374 Video pixel networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*,
375 pp. 1771–1779. JMLR. org, 2017.
- 376 Kingma, D. P. and Welling, M. Auto-encoding variational bayes. In *International Conference on Learning
377 Representations*, 2013.
- 378 Krizhevsky, A. Learning multiple layers of features from tiny images. 2009.
- 379 Kunin, D., Bloom, J. M., Goeva, A., and Seed, C. Loss landscapes of regularized linear autoencoders. *ArXiv*,
380 abs/1901.08168, 2019.
- 381 Liu, Z., Luo, P., Wang, X., and Tang, X. Deep learning face attributes in the wild. *2015 IEEE International
382 Conference on Computer Vision (ICCV)*, pp. 3730–3738, 2014.
- 383 Menick, J. and Kalchbrenner, N. Generating high fidelity images with subscale pixel networks and multidimen-
384 sional upscaling. *arXiv preprint arXiv:1812.01608*, 2018.
- 385 Oord, A. v. d., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A., and
386 Kavukcuoglu, K. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016a.
- 387 Oord, A. v. d., Kalchbrenner, N., and Kavukcuoglu, K. Pixel recurrent neural networks. *arXiv preprint
388 arXiv:1601.06759*, 2016b.
- 389 Papamakarios, G., Pavlakou, T., and Murray, I. Masked autoregressive flow for density estimation. In *Advances
390 in Neural Information Processing Systems*, pp. 2338–2347, 2017.
- 391 Pfau, D., Petersen, S., Agarwal, A., Barrett, D. G. T., and Stachenfeld, K. L. Spectral inference networks:
392 Unifying deep and spectral learning. In *ICLR*, 2019.
- 393 Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. Language models are unsupervised
394 multitask learners. *OpenAI Blog*, 1(8):9, 2019.
- 395 Ramachandran, P., Paine, T. L., Khorrami, P., Babaeizadeh, M., Chang, S., Zhang, Y., Hasegawa-Johnson,
396 M. A., Campbell, R. H., and Huang, T. S. Fast generation for convolutional autoregressive models. *ArXiv*,
397 abs/1704.06001, 2017.
- 398 Razavi, A., van den Oord, A., and Vinyals, O. Generating diverse high-fidelity images with vq-vae-2. In
399 *Advances in Neural Information Processing Systems*, pp. 14837–14847, 2019.
- 400 Rippel, O., Gelbart, M. A., and Adams, R. P. Learning ordered representations with nested dropout. In *ICML*,
401 2014.

- 402 Salimans, T., Karpathy, A., Chen, X., and Kingma, D. P. Pixelcnn++: Improving the pixelcnn with discretized
403 logistic mixture likelihood and other modifications. *arXiv preprint arXiv:1701.05517*, 2017.
- 404 Scholz, M. and Vigário, R. Nonlinear pca: a new hierarchical approach. In *ESANN*, 2002.
- 405 van den Oord, A., Kalchbrenner, N., Espeholt, L., Kavukcuoglu, K., Vinyals, O., and Graves, A. Conditional
406 image generation with pixelcnn decoders. In *NIPS*, 2016.
- 407 van den Oord, A., Vinyals, O., and Kavukcuoglu, K. Neural discrete representation learning. In *NIPS*, 2017.
- 408 van den Oord, A., Li, Y., Babuschkin, I., Simonyan, K., Vinyals, O., Kavukcuoglu, K., van den Driessche,
409 G., Lockhart, E., Cobo, L. C., Stimberg, F., Casagrande, N., Grewe, D., Noury, S., Dieleman, S., Elsen, E.,
410 Kalchbrenner, N., Zen, H., Graves, A., King, H., Walters, T., Belov, D., and Hassabis, D. Parallel wavenet:
411 Fast high-fidelity speech synthesis. *ArXiv*, abs/1711.10433, 2018.
- 412 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I.
413 Attention is all you need. In *Advances in neural information processing systems*, pp. 5998–6008, 2017.
- 414 Veaux, C., Yamagishi, J., and Macdonald, K. Cstr vctk corpus: English multi-speaker corpus for cstr voice
415 cloning toolkit. 2017.