
Computational Tradeoffs in Image Synthesis: Diffusion, Masked-Token, and Next-Token Prediction

Maciej Kilian^π Varun Jampani^π Luke Zettlemoyer^μ
^πStabilityAI ^μUniversity of Washington

Abstract

Nearly every recent image synthesis approach, including diffusion, masked-token prediction, and next-token prediction, uses a Transformer network architecture. Despite this common backbone, there has been no direct, compute controlled comparison of how these approaches affect performance and efficiency. We analyze the scalability of each approach through the lens of compute budget measured in FLOPs. We find that token prediction methods, led by next-token prediction, significantly outperform diffusion on prompt following. On image quality, while next-token prediction initially performs better, scaling trends suggest it is eventually matched by diffusion. We compare the inference compute efficiency of each approach and find that next token prediction is by far the most efficient. Based on our findings we recommend diffusion for applications targeting image quality and low latency; and next-token prediction when prompt following or throughput is more important.

1 Introduction

Following the work of Peebles and Xie [2023], deep image synthesis, including diffusion [Sohl-Dickstein et al., 2015, Song and Ermon, 2019, Song et al., 2020, Ho et al., 2020, Rombach et al., 2022, Esser et al., 2024], masked-token prediction [Chang et al., 2022, 2023, Villegas et al., 2022, Yu et al., 2024], and next-token prediction [Gafni et al., 2022, Yu et al., 2022b, Esser et al., 2021], are all build on a common Transformer architecture [Vaswani et al., 2023]. Although these approaches are all known to scale well with compute and data, there has been relatively little controlled comparisons of their relative training and inference efficiency. Comparing these latent image synthesis approaches is challenging since the objectives they optimize often have different requirements which limit the set of applicable modules for each approach and influence their optimal configurations. For example, next-token prediction requires discrete input data which makes it unfit for continuous latent space regularization advances. In fact, latent image synthesis will be strongly influenced by the state of autoencoding research, often in an unbalanced way. Examples of this can be found in Section 2.

In this paper, we measure the computational tradeoffs between popular transformer-based latent image synthesis approaches - diffusion, masked-token prediction, and next-token prediction. We investigate the impact of the autoencoder, which encodes the latent space, on generative results and train a large grid of models with the different approaches, model sizes, and dataset sizes. Samples from some of our most capable models can be found in Figure 1. Our findings indicate that (i) at smaller compute budgets, next-token prediction yields the best image quality but scaling trends suggest it is eventually matched by diffusion. (ii) Token-based approaches achieve superior controllability. (iii) The quality of the autoencoder impacts the FID more than the CLIP score of diffusion models trained on its latent space. (iv) We find preliminary evidence for improved diffusion training practices. Based on our findings, we recommend diffusion models for applications targeting low latency and high image quality; and next-token prediction for applications where prompt following and throughput are priorities.



Figure 1: **Images generated using our best models.** Top row is from a next-token prediction model, bottom row is from a diffusion model. Both models are XL size and trained for 500k steps.

2 Related Work

Scaling transformer-based generative models. Scaling compute budgets for transformer based generative models is a predictable method for improving performance. Kaplan et al. [2020], Hoffmann et al. [2022], Clark et al. [2022] showed that for text, final training loss can be accurately predicted as a power law of training compute which depends on model size and dataset size. Following those practices many capable text generation models were trained [Touvron et al., 2023, Brown et al., 2020, Rae et al., 2022]. Similar results have been found for vision [Zhai et al., 2022, Alabdulmohsin et al., 2024, Esser et al., 2024, Dehghani et al., 2023] and even mixed modal data [Aghajanyan et al., 2023]. We follow these intuitions and analyze image synthesis performance as a function of compute budget.

Latent generative modeling. Training latent generative vision models has emerged as an efficient alternative to the computationally intensive modeling of high-dimensional pixel space. Studies have demonstrated the advantages of imposing specific structural regularizations within the latent space for enhancing the performance of various generative models. For instance, Rombach et al. [2022] observed that latent diffusion models operating in VAE-style Kingma and Welling [2022] latent spaces, when regularized towards a standard Gaussian structure, outperform models trained with alternative regularization techniques. Yu et al. [2024], Mentzer et al. [2023], Yu et al. [2022a] have shown that simplifying vector quantization methods can mitigate common issues such as poor codebook utilization and enhancing the transfer between autoencoder reconstruction quality and downstream generative model performance for token-based approaches. Tian et al. [2024] demonstrated that employing hierarchical next-scale latents enables transformers using next token prediction to leverage their in-context learning capabilities more effectively, significantly improving performance. Jin et al. [2024] use image latents dynamically sized based on their information content which allows generative models to allocate more computational resources to complex samples, as these will contain more tokens. We minimize potential bias coming from autoencoding asymmetries by studying the impact of the autoencoder on the generative model trained on top of it.

3 Background

Autoencoding To train latent generative models, we establish an encoder-decoder pair $(\mathcal{E}, \mathcal{D})$. For an image $x \in \mathbb{R}^{H \times W \times 3}$, the encoder maps x to a latent representation $z = \mathcal{E}(x)$, where $z \in \mathbb{R}^{H/f \times W/f \times c}$ and $f = 2^{m \in \mathbb{N}}$ represents the factor of dimensionality reduction. The decoder then reconstructs $\hat{x} = \mathcal{D}(z)$, aiming for high perceptual similarity to x , effectively making z a perceptually compressed representation of the input. To avoid high-variance latent spaces and ensure structured representations, we employ regularization methods classified into two main types: discrete and continuous. The regularization function \mathbf{q} for a continuous regularizer maps $\mathbb{R}^d \rightarrow \mathbb{R}^d$, while a discrete regularizer maps $\mathbf{q} : \mathbb{R}^d \rightarrow \{0, 1, 2, \dots, N\}$, making the latent space finite. In the following subsections we use $z \in \mathbb{R}^{s \times d}$ to denote the flattened representation output by the encoder \mathcal{E} . $p(z)$ is the latent data distribution we are interested in estimating using our generative models. We recover images by inputting sampled latents into the corresponding decoder \mathcal{D} .

Next token prediction In the context of sequences of discrete tokens represented as $z \in \{0, 1, 2, \dots, N\}^s$, we employ the chain rule of conditional probability to decompose the target distribution into a product of conditional distributions which are tractable since the range of z_i is finite. To model this distribution, we use a neural network f , parameterized by weights θ . The

parameters are optimized by minimizing the negative log-likelihood \mathcal{L}_{NT} .

$$p(z) = \prod_{i=1}^n p(z_i | z_{i-1}, \dots, z_1) \quad \mathcal{L}_{NT} = \mathbb{E}_i[-\log p(z_i | z_{<i}; \theta)] \quad (1)$$

Sampling from our learned distribution begins with an empty sequence (in practice, a "start of text" token is sampled with 1.0 probability). We then sample the first token unconditionally and append it to our sequence. The process continues by iteratively evaluating the conditionals and sampling from them, with each step increasing the sequence length by one.

Masked token prediction Masked token prediction is a form of iterative denoising and can be viewed as a discrete diffusion process. In this process, tokens progressively transition to an absorbing [MASK] state according to a probability defined by a noise schedule $\gamma(t) \in (0, 1]$ where $t \sim \mathcal{U}(0, 1)$. This transition can also be mathematically expressed as a product of conditionals, except in a perturbed order σ , and implemented as a neural network. Here, $\sigma(i)$ is a surjective function mapping $[0, N] \mapsto [0, N]$. We follow Chang et al. [2022, 2023] where $\sigma(i) = \sigma(i, t)$ such that $p(\sigma(i, t) < j) = \gamma(t)$ meaning the likelihood a token can be attended to is independent of position. In this formulation, we utilize a truncated arccos distribution for our noise schedule: $\gamma(t) = \frac{2}{\pi}(1 - t^2)^{-\frac{1}{2}}$. To apply this method, we generate a mask tensor $M \in \{0, 1\}^s$ by sampling $t \sim \mathcal{U}(0, 1)$ and $m_i \sim \text{Bernoulli}(\gamma(t))$. The tensor M is applied elementwise to the latents, replacing z_i with the [MASK] token if $m_i = 1$; otherwise, z_i remains unchanged. Denote the resultant noised sequence as z_M . The network is then trained to minimize the masked token loss \mathcal{L}_{MT} .

$$p(z) = \prod_{i, \sigma(i)=j}^n p(z_j | z_{\sigma(i) < j}) \quad \mathcal{L}_{MT} = \mathbb{E}_{i, m_i=1}[-\log p(z_i | z_{\overline{M}}; \theta)] \quad (2)$$

Sampling from the distribution starts with a fully masked sequence and iterates through a discretized noise schedule $t_i = i/N$ over N desired steps. At each step, the model estimates $p(z | z_{\overline{M}})$ for sampling, followed by re-noising using $\gamma(t_{i+1})$. This iterative re-noising and sampling process is repeated N times to yield the final sample.

Diffusion We adopt the flow matching framework outlined by Lipman et al. [2023], focusing on models that map samples from a noise distribution p_1 to a data distribution p_0 using continuous trajectories governed by an ordinary differential equation (ODE). Furthermore, we enforce straight paths between the terminal distributions (by setting $\alpha_t = 1 - t$ and $\beta_t = t$) since this has been shown to perform well at scale [Esser et al., 2024].

$$d\phi_t(x) = v_t(\phi_t(x)) dt \quad \phi_0(x) = x \quad z_t = \alpha_t x_0 + \beta_t \epsilon, \quad \epsilon \sim \mathcal{N}(0, 1) \quad (3)$$

Here, $v_t : [0, 1] \times \mathbb{R}^d \mapsto \mathbb{R}^d$ represents a time-dependent vector field, which we aim to parameterize using a neural network θ and $\phi_t : [0, 1] \times \mathbb{R}^d \mapsto \mathbb{R}^d$ is the flow. To optimize our weights, we regress the vector field u_t , which generates our paths z_t by employing conditional flow matching which we reformulate as a noise-prediction objective \mathcal{L}_{CFM} . Sampling is performed by using an ODE solver to solve Equation 3 in reverse time, utilizing our trained neural network $v_\theta(z, t)$.

$$\mathcal{L}_{CFM}(x_0) = \mathbb{E}_{t \sim \mathcal{U}(0, 1), \epsilon \sim \mathcal{N}(0, 1)} \left[\frac{1}{(1 - t)^2} \|\epsilon_\theta(z_t, t) - \epsilon\|^2 \right]$$

4 Experimental Setup

Data We train both the autoencoders and the generative models on a large web dataset of image and text pairs at 256x256 resolution. For the conditioning we use the pooled text embedding of the OpenCLIP bigG/14 model from Cherti et al. [2022]. Once the autoencoders are trained we pre-encode the entire dataset with them for improved training speed.

Evaluation metrics Since we are in the infinite data regime, we look at the final train loss and do not compare across objectives since the losses represent different quantities. We also look at CLIP score [Radford et al., 2021, Hessel et al., 2022] and FID computed on CLIP features [Sauer et al., 2021] based on the decoded samples $\hat{x} = \mathcal{D}(z)$.

Regularizer	Latent space capacity	rFID (\downarrow)	Model size	Layers N	Hidden size d	Heads
KL	16 channels	1.060	S	12	768	12
KL	8 channels	1.560	M	24	1024	16
KL _{es}	8 channels	2.856	L	24	1536	16
KL	4 channels	2.410	XL	32	2304	32
LFQ	16384 vocabulary	2.784				

Table 1: **Autoencoders.** Reconstruction metrics for differently regularized and trained autoencoders (downsampling $f = 8$). "es" is early stopping to match the LFQ autoencoder.

Table 2: **Transformer configurations.** Base transformer hyperparameters for models we train. Common across all approaches

Autoencoding We study well-established autoencoder configurations that have proven effective without special handling for each data type. We adhere to the training and architectural guidelines provided by Rombach et al. [2022]. Each autoencoder is trained with a downsampling factor $f = 8$, reducing 256×256 images to a 32×32 grid of latents. For continuous variants, $\mathbf{q}(z)$ implements a KL penalty aiming towards the standard normal distribution [Kingma and Welling, 2022, Rombach et al., 2022], while for discrete variants, we utilize lookup-free quantization (LFQ) [Yu et al., 2024]. Further details on the selection of discrete regularizers are available in Appendix A. To circumvent potential challenges associated with large vocabulary sizes, as highlighted by Yu et al. [2024], our LFQ-regularized autoencoder is trained with a vocabulary size of 16384 [Esser et al., 2021]. Assessing the comparability of autoencoders is difficult since there are many variables of interest such as the (1) information capacity of the latent space; (2) compute used to train the autoencoder; (3) reconstruction quality achieved by the autoencoder. To explore the influence of these factors on the performance of generative models, we train a set of autoencoders similar to those in Esser et al. [2024], which exhibit a range of information capacities and reconstruction qualities. Additionally, we experiment with targeting specific reconstruction qualities, irrespective of other factors, by training a KL-regularized autoencoder with early stopping to match the reconstruction quality of our discrete autoencoder within a certain threshold ϵ ¹. Table 1 provides detailed information about the autoencoders.

Autoencoder ablation. We train an L-size diffusion model on top of the latent space of each continuous autoencoder. We then evaluate the models using the metrics described in Section 4 and plot them against the number of training steps. Results are shown in Figure 2. We find that the autoencoder’s reconstruction quality has a consistently significant impact on the FID score, while its effect on the CLIP score diminishes with larger dataset sizes, where the models tend to yield similar results. This trend likely emerges because improvements in autoencoder quality enhance perceptual reconstruction metrics similar to FID, rather than affecting language or semantic capabilities. Upon examining the number of channels in the autoencoders, our findings concur with those reported by Esser et al. [2024], indicating that leveraging larger and better latent spaces requires more compute and model capacity. Additionally, the model trained on our early-stopped autoencoder’s latent space performed significantly worse than the 4-channel autoencoder, which achieves similar reconstruction quality. This confirms the importance of latent space structure for overall performance. Building on these insights, we have chosen to use the 4-channel autoencoder for our main diffusion experiments. This model most closely matches the latent space capacity and reconstruction quality of our discrete autoencoder, while also ensuring that the latent structure is adequately developed to support the diffusion model trained on it. Although more advanced autoencoders have been developed—such as those featuring increased channel counts or expanded codebook sizes—our primary focus in this study is to maintain comparability across objectives.

4.1 Network Architecture

Backbone. We opt for the transformer architecture as our primary network backbone, recognizing its capability to scale effectively with computational resources and its status as the state-of-the-art (SOTA) across all evaluated approaches. Configuring a transformer involves many decisions, such as choosing normalization methods, feed-forward layer configurations, positional embedding schemes,

¹Discrete autoencoders typically have worse reconstruction qualities since the information bottleneck is tighter. This can be shown by comparing $\log(\text{codebook size})$ to $\text{num_channels} * \text{sizeof}(\text{dtype})$ for common values of these quantities. In our case we needed to stop at 75k steps vs. 1M for the discrete autoencoder.

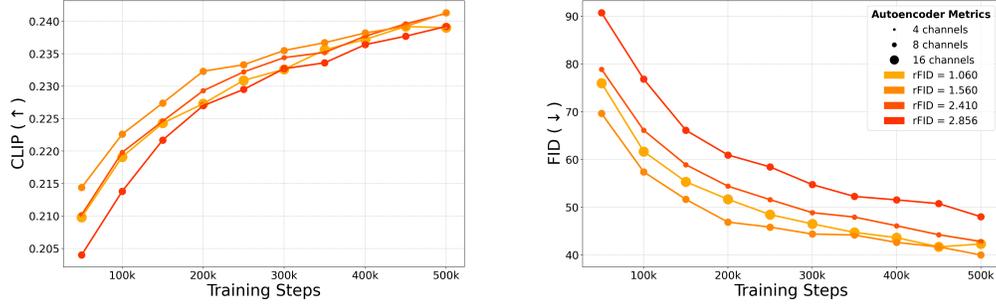


Figure 2: **Impact of autoencoder quality on diffusion models.** We train L-size diffusion models on our set of continuous latent space autoencoders. The choice of autoencoder has more impact on FID than CLIP score. Effectively using a larger latent space requires more compute and model capacity.

Model	n-parameters (%)	Forward TFLOPs
DiT-S	131.13 M (97.2%)	0.2133
DiT-M	459.19 M (98.7%)	0.7234
DiT-L	1031.67 M (98.8%)	1.5485
DiT-XL	3083.69 M (99.2%)	4.4901
NT/MT-S	153.73 M (82.9%)	0.2261
NT/MT-M	494.56 M (92.9%)	0.6631
NT/MT-L	1072.14 M (95.1%)	1.3421
NT/MT-XL	3137.29 M (97.5%)	3.7166

Table 3: **Model forward pass costs.** Number of parameters and FLOPs used in each forward pass for all models we trained. DiT - Diffusion Transformer; NT - Next-token; MT - Masked-token

Objective	Conditioning	FID	CLIP
NT	adaLNzero	83.052	0.2213
NT	in context	88.176	0.2041
NT	cross attention	92.852	0.2062
MT	adaLNzero	97.021	0.2164
MT	in context	100.646	0.1925
MT	cross attention	103.221	0.1960

Table 4: **Conditioning method ablation.** Results for different objectives and conditioning methods. adaLNzero conditioning is used for the remainder of experiments. NT - Next-token; MT - Masked-token.

conditioning methods, and initialization strategies. Given the prohibitive cost of exploring all possible hyperparameters, we adhere to established practices in recent studies.

Design differences. For approaches utilizing discrete representations, we primarily follow the configurations used in the LLaMa model [Touvron et al., 2023], incorporating SwiGLU feed-forward layers with an expansion ratio of $\frac{2}{3} \cdot 4$ and rotary positional embeddings [Su et al., 2023]. An exception is made for masked token prediction, where learned positional embeddings are preferred to address positional ambiguities that degrade performance near the center of the image. For continuous representation approaches, we align with diffusion transformers [Peebles and Xie, 2023], employing GELU feed-forward layers with an expansion ratio of 4 and learned positional embeddings. All models use QK-normalization [Dehghani et al., 2023] for better training stability.

Conditioning ablation. We choose to ablate the conditioning method, as it significantly impacts the computational cost of model operations. Adaptive layer normalization (AdaLN) [Perez et al., 2017] has shown promise in latent image synthesis for both continuous [Peebles and Xie, 2023] and discrete [Tian et al., 2024] settings. To validate this choice in the discrete context, we conduct small-scale ablations on S-size models, comparing AdaLNzero [Peebles and Xie, 2023], with two other common conditioning methods: prepending a projected embedding in the context of the transformer and cross-attention. The outcomes of these ablations are presented in Table 4, informing our choice of conditioning method for subsequent experiments.

Compute cost. To assess the computational cost of each model, we first standardize a set of hyperparameters across all transformers, detailed in Table 2. We then calculate the forward pass FLOPs for a single sample (a sequence of 1024 embeddings) for each approach and model size, and present them in Table 3. Assuming the backward pass is twice the cost of the forward pass, we compute the training FLOPs for each model as $(1 + 2) \times (\text{forward FLOPs}) \times D$, where D represents the total number of training samples.

4.2 Training

Each approach also has associated training hyperparameters which past work has found to work well and for the same reasons as stated in 4.1 we follow them.

Optimization and conditioning. For diffusion experiments we follow Esser et al. [2024] and use a constant learning rate schedule with a maximum value of 1^{-4} . For next and masked token prediction we use a cosine decay learning rate with a maximum value of 3^{-3} which decays down to 3^{-5} . All models have a linear learning rate warmup lasting 1000 steps up to the maximum value. We use the AdamW [Loshchilov and Hutter, 2019] optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.95$, decay=0.01, and epsilon=1e-15 for improved transformer training stability [Wortsman et al., 2023]. All models are trained at bf16-mixed precision [Chen et al., 2019]. We intend to use classifier free guidance (CFG) [Ho and Salimans, 2022] during sampling so we randomly drop conditioning 10% of the time during training. Since its inexpensive and does not influence training, for all models, we store a copy of the model weights which gets updated every 100 training batches with an exponential moving average (EMA) using a decay factor of 0.99 and during evaluation we evaluate both sets of weights.

Training steps. For each objective and model size we scale to at least 250k training steps with a batch size of 512. For diffusion we decide to go up to 500k steps since constant learning rate schedule allows more flexibility with dataset size². Occasionally we train models for longer to attempt to illustrate convergence or crossing points.

4.3 Sampling

Classifier free guidance. Ho and Salimans [2022] introduced it in diffusion models as an elegant way of trading off diversity for fidelity and has been demonstrated to improve results for all approaches we consider in this study [Chang et al., 2023, Gafni et al., 2022, Ho and Salimans, 2022]. We use it here in the form

$$x_g = (1 + w)x_c - wx_u \tag{4}$$

where w is the guidance scale. For diffusion x will be the position in the denoising trajectory and for token based methods x is the logit distribution at a given timestep.

Hyperparameters. For our diffusion models we follow Esser et al. [2024] and use 50 sampling steps with a CFG scale of 5. Since the conditioning and input data is slightly different we also perform a small sweep around those parameters to confirm they are still optimal. For the token based models we could not find good resources on reasonable sampling hyperparameters so we perform small sweeps for S-size models to find the best configurations and verify the robustness of those values for larger models. Common between them, we use nucleus sampling [Holtzman et al., 2020] with a top-p value of 0.9 and a temperature of 1.0. For next token prediction and masked token prediction we use CFG scales 8 and 5 respectively. For masked token prediction we perform 10 sampling steps.

5 Results

5.1 Training tradeoffs

For all models, we measure our evaluation metrics every 50k steps of training and plot them in log scale against the log of training compute. Figure 3 presents this for FID and CLIP score. There we can see that for FID, next token prediction starts out more compute efficient but scaling trends suggest that its eventually matched by diffusion. When looking at CLIP score we see that token prediction is significantly better than diffusion, implying the models generate images that follow the input prompt better. This could be a feature of using more compressed latent spaces which is supported by Figure 2 where the 4 channel continuous autoencoder outperforms both the 8 and 16 channel autoencoder on CLIP score near the end of training. This is also supported in Figure 7 with interpretable features like human faces emerging sooner in the token based methods. Extending a finding from Mei et al. [2024], we observe that, for all approaches studied, smaller models trained for longer often surpass larger models. In Figure 4 we show the final training loss of each model against training compute to show that it follow similar scaling trends to what has been shown in past work on scaling deep neural networks, briefly described in Section 2. Samples from the most capable XL sized next-token prediction and diffusion models can be found in Figure 1.

²With a decaying learning rate, each dataset size we want to study requires a separate run from scratch whereas for constant learning rate schedules you can simply continue from a past checkpoint

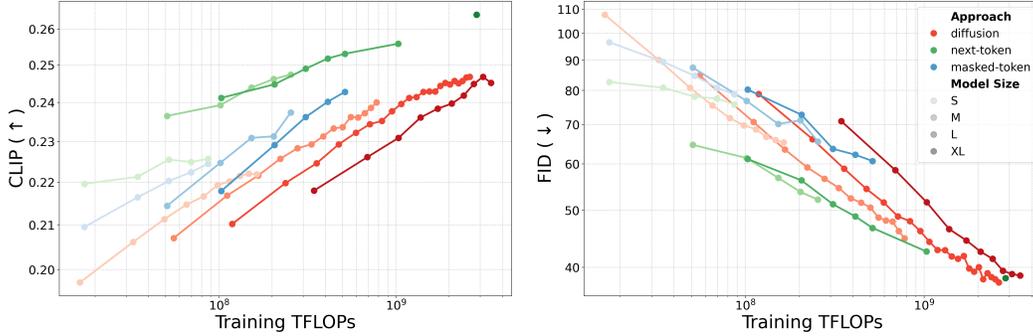


Figure 3: **Training compute efficiency on perceptual metrics.** Performance on CLIP and FID scores for various models and dataset sizes across different image synthesis approaches. On FID, next-token prediction is initially the most compute-efficient but scaling trends suggest it is eventually matched by diffusion. Token-based methods significantly outperform diffusion in CLIP score. Both axes are in log scale.

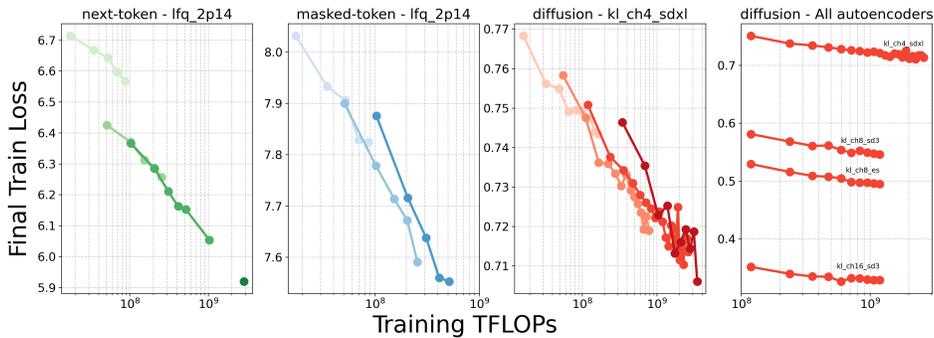


Figure 4: **Training compute efficiency on final loss.** All objectives follow predictable scaling trends. Right plot shows the difference in loss scale between diffusion models trained on top of different autoencoders. FLOPs axis is in log scale.

5.2 Inference tradeoffs

Inference cost. We evaluated all models trained for 250k steps to understand the impact of inference FLOPs on perceptual metrics. To adjust the number of inference FLOPs for a single model, we varied the number of sampling steps, applicable only to iterative denoising methods like masked token prediction and diffusion. As shown in Figure 5, next-token prediction demonstrates far greater inference compute efficiency compared to other objectives. This efficiency arises because when using key-value caching, sampling N tokens autoregressively uses the same amount of FLOPs as forwarding those N tokens in parallel once. However, for iterative denoising methods, this value is multiplied by the number of sampling steps. Interestingly, despite being trained for iterative denoising, the number of steps in masked token prediction appears to have minimal impact on sample quality.

Sampling latency and throughput. While next-token prediction requires much less compute per sample, the autoregressive dependency of each token causes it to be data bound when few queries are being processed in parallel which results in high latency. Conversely, bidirectional denoising approaches utilize a more parallel sampling process which, despite its high cost, facilitates low latency especially in low-volume settings with models that fit on local devices [Chang et al., 2022]. For high-volume sampling, where throughput becomes more important, such as serving many users via an API, next token prediction could use a batching algorithm to maximize GPU utilization by choosing batch sizes inversely proportional to sequence lengths. The effectiveness of this method is ensured by the fact that, for next-token prediction image synthesis, all responses are the same length so you can easily plan your batches ahead. This way, for high-volume sampling, next-token prediction would enjoy the same benefits over the other approaches as presented in the cost section above but for sample throughput.

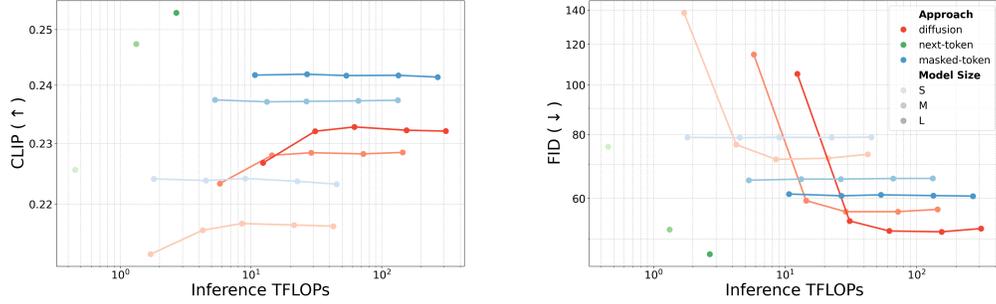


Figure 5: **Inference compute efficiency on perceptual metrics.** Diffusion and masked token prediction evaluated at 4, 10, 20, 50, and 100 sampling steps. Next token prediction is 1 forward pass factorized over each token individually. Masked token prediction isn't influenced by the number of sampling steps very much. Next token prediction is the most compute efficient. Both axes are in log scale.

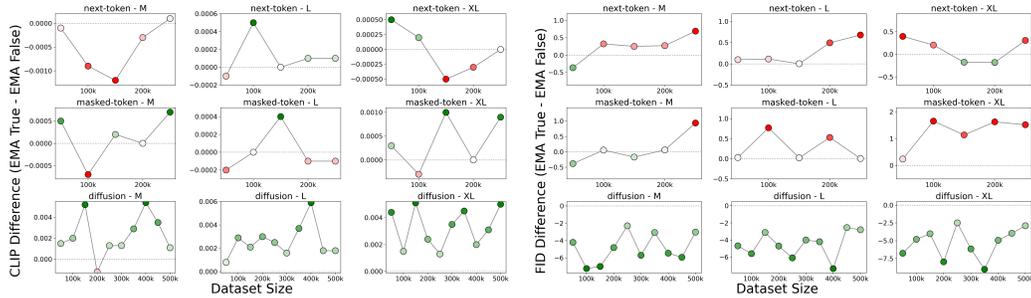


Figure 6: **Impact of EMA.** EMA significantly improves FID for diffusion models but hurts token based approaches. On CLIP score the effect on diffusion models stays consistent however for token based methods the influence is negligible.

5.3 EMA ablations

Among the various training practices distinguishing these methods, the use of an exponential moving average (EMA) on the model weights stands out. In the diffusion literature [Karras et al., 2024, 2022, Peebles and Xie, 2023, Esser et al., 2024] EMA is an essential component of the training pipeline. In contrast, this practice has not received equivalent attention in other approaches. The differential impact of EMA is evident in Figure 6. For token-based approaches, the influence of EMA is either negligible or, in some cases, harmful, whereas for diffusion models, it is beneficial almost universally. We hypothesize that the impact of EMA may be linked to the learning rate schedule, where decaying schedules similarly minimize weight variation towards the end of training. To test this hypothesis, we conducted an ablation study on an M-sized next-token prediction and diffusion model trained over 250k steps. Our findings verify our hypothesis that EMA enhances performance under a constant learning rate schedule; however, it does not exceed the improvements seen with a cosine decay learning rate schedule. This implies that future diffusion models should consider substituting the EMA for a cosine decay learning rate schedule if they are willing to pay the cost of decreased training length flexibility. Results from this ablation study are presented in Table 5.

5.4 Limitations

Our analysis has several limitations which result from resource limitations and project scope. We only investigate pretraining whereas most production systems utilize a progression of pretraining, finetuning, and distillation stages. We do not investigate high resolution images. We only measure loss and perceptual metrics and leave out an analysis of utility for potential downstream tasks. There are many others approaches that we leave out such as other discrete diffusion approaches [Austin et al., 2023, Pernias et al., 2023], causally masked token prediction [Aghajanyan et al., 2022], and many more. We choose most hyperparameters by following past work instead of exhaustively sweeping

Objective	LR schedule	EMA	FID	CLIP
Next-token	constant	✗	81.976	0.2208
Next-token	constant	✓	79.571	0.2230
Next-token	cosine	✗	75.715	0.2256
Next-token	cosine	✓	76.404	0.2257
Diffusion	constant	✗	74.087	0.2153
Diffusion	constant	✓	71.789	0.2166
Diffusion	cosine	✗	69.284	0.2195
Diffusion	cosine	✓	69.468	0.2192

Table 5: **EMA and learning rate schedules.** EMA on model weights improves results under a constant learning rate schedule but does not exceed the gains from using a cosine decay schedule.

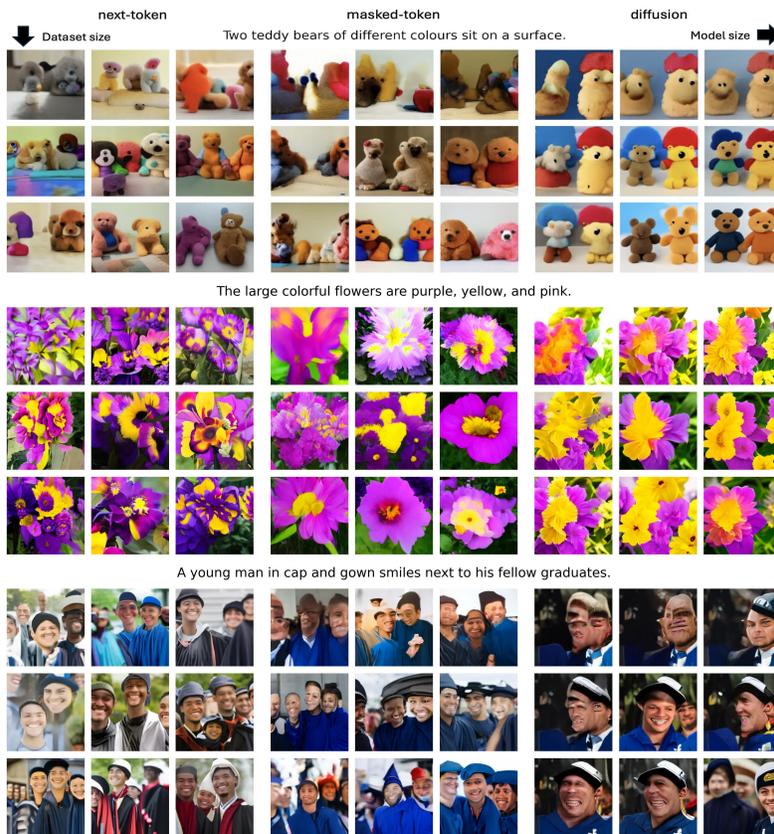


Figure 7: **Increasing training compute improves sample quality for all approaches.** For each approach and prompt we sample an image with all combinations of S, M, L model sizes and 50k, 150k, 250k dataset sizes. Going down or right in the 3x3 increases dataset and model size respectively.

to find the best configurations. And finally, we do not compare approaches using the best possible autoencoders.

6 Conclusion

We conduct a compute-controlled analysis comparing transformer-based diffusion, next-token prediction, and masked-token prediction latent image synthesis models. Our findings indicate that token based methods, led by next-token prediction, achieve superior CLIP scores, indicating greater controllability. In terms of FID, and therefore image quality, while next-token prediction is much better at low training compute scales, scaling trends suggest it is eventually matched by diffusion.

We find that next token prediction has, by far, the best inference compute efficiency but this comes at the cost of high latency in low data intensity settings. Based on our findings recommend diffusion models when image quality and low latency is important; and next-token prediction for better prompt following and throughput.

References

- Armen Aghajanyan, Bernie Huang, Candace Ross, Vladimir Karpukhin, Hu Xu, Naman Goyal, Dmytro Okhonko, Mandar Joshi, Gargi Ghosh, Mike Lewis, and Luke Zettlemoyer. Cm3: A causal masked multimodal model of the internet, 2022.
- Armen Aghajanyan, Lili Yu, Alexis Conneau, Wei-Ning Hsu, Karen Hambardzumyan, Susan Zhang, Stephen Roller, Naman Goyal, Omer Levy, and Luke Zettlemoyer. Scaling laws for generative mixed-modal language models, 2023.
- Ibrahim Alabdulmohsin, Xiaohua Zhai, Alexander Kolesnikov, and Lucas Beyer. Getting vit in shape: Scaling laws for compute-optimal model design, 2024.
- Jacob Austin, Daniel D. Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. Structured denoising diffusion models in discrete state-spaces, 2023.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
- Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T. Freeman. Maskgit: Masked generative image transformer, 2022.
- Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T. Freeman, Michael Rubinstein, Yuanzhen Li, and Dilip Krishnan. Muse: Text-to-image generation via masked generative transformers, 2023.
- Dehao Chen, Chiachen Chou, Yuanzhong Xu, and Jonathan Hseu. Bfloat16: The secret to high performance on cloud tpus, 2019. URL <https://cloud.google.com/blog/products/ai-machine-learning/bfloat16-the-secret-to-high-performance-on-cloud-tpus?hl=en>.
- Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning, 2022.
- Aidan Clark, Diego de las Casas, Aurelia Guy, Arthur Mensch, Michela Paganini, Jordan Hoffmann, Bogdan Damoc, Blake Hechtman, Trevor Cai, Sebastian Borgeaud, George van den Driessche, Eliza Rutherford, Tom Hennigan, Matthew Johnson, Katie Millican, Albin Cassirer, Chris Jones, Elena Buchatskaya, David Budden, Laurent Sifre, Simon Osindero, Oriol Vinyals, Jack Rae, Erich Elsen, Koray Kavukcuoglu, and Karen Simonyan. Unified scaling laws for routed language models, 2022.
- Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, Rodolphe Jenatton, Lucas Beyer, Michael Tschannen, Anurag Arnab, Xiao Wang, Carlos Riquelme, Matthias Minderer, Joan Puigcerver, Utku Evci, Manoj Kumar, Sjoerd van Steenkiste, Gamaleldin F. Elsayed, Aravindh Mahendran, Fisher Yu, Avital Oliver, Fantine Huot, Jasmijn Bastings, Mark Patrick Collier, Alexey Gritsenko, Vighnesh Birodkar, Cristina Vasconcelos, Yi Tay, Thomas Mensink, Alexander Kolesnikov, Filip Pavetić, Dustin Tran, Thomas Kipf, Mario Lučić, Xiaohua Zhai, Daniel Keysers, Jeremiah Harmsen, and Neil Houlsby. Scaling vision transformers to 22 billion parameters, 2023.
- Patrick Esser, Robin Rombach, and Björn Ommer. Taming transformers for high-resolution image synthesis, 2021.

Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis, 2024.

Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors, 2022.

Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning, 2022.

Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance, 2022.

Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.

Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models, 2022.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration, 2020.

Yang Jin, Kun Xu, Kun Xu, Liwei Chen, Chao Liao, Jianchao Tan, Quzhe Huang, Bin Chen, Chenyi Lei, An Liu, Chengru Song, Xiaoqiang Lei, Di Zhang, Wenwu Ou, Kun Gai, and Yadong Mu. Unified language-vision pretraining in llm with dynamic discrete visual tokenization, 2024.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models, 2020.

Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models, 2022.

Tero Karras, Miika Aittala, Jaakko Lehtinen, Janne Hellsten, Timo Aila, and Samuli Laine. Analyzing and improving the training dynamics of diffusion models, 2024.

Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2022.

Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling, 2023.

Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.

Kangfu Mei, Zhengzhong Tu, Mauricio Delbracio, Hossein Talebi, Vishal M. Patel, and Peyman Milanfar. Bigger is not always better: Scaling properties of latent diffusion models, 2024.

Fabian Mentzer, David Minnen, Eirikur Agustsson, and Michael Tschannen. Finite scalar quantization: Vq-vae made simple, 2023.

William Peebles and Saining Xie. Scalable diffusion models with transformers, 2023.

Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer, 2017.

Pablo Pernias, Dominic Rampas, Mats L. Richter, Christopher J. Pal, and Marc Aubreville. Wuerstchen: An efficient architecture for large-scale text-to-image diffusion models, 2023.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021.

- Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu, Erich Elsen, Siddhant Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsimpoukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d’Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew Johnson, Blake Hechtman, Laura Weidinger, Iason Gabriel, William Isaac, Ed Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub, Jeff Stanway, Lorraine Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. Scaling language models: Methods, analysis and insights from training gopher, 2022.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2022.
- Axel Sauer, Kashyap Chitta, Jens Müller, and Andreas Geiger. Projected gans converge faster, 2021.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, 2023.
- Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive modeling: Scalable image generation via next-scale prediction, 2024.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023.
- Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning, 2018.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023.
- Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang, Mohammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. Phenaki: Variable length video generation from open domain textual description, 2022.
- Mitchell Wortsman, Peter J. Liu, Lechao Xiao, Katie Everett, Alex Alemi, Ben Adlam, John D. Co-Reyes, Izzeddin Gur, Abhishek Kumar, Roman Novak, Jeffrey Pennington, Jascha Sohl-dickstein, Kelvin Xu, Jaehoon Lee, Justin Gilmer, and Simon Kornblith. Small-scale proxies for large-scale transformer training instabilities, 2023.
- Jiahui Yu, Xin Li, Jing Yu Koh, Han Zhang, Ruoming Pang, James Qin, Alexander Ku, Yuanzhong Xu, Jason Baldridge, and Yonghui Wu. Vector-quantized image modeling with improved vqgan, 2022a.

Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling autoregressive models for content-rich text-to-image generation, 2022b.

Lijun Yu, José Lezama, Nitesh B. Gundavarapu, Luca Versari, Kihyuk Sohn, David Minnen, Yong Cheng, Vighnesh Birodkar, Agrim Gupta, Xiuye Gu, Alexander G. Hauptmann, Boqing Gong, Ming-Hsuan Yang, Irfan Essa, David A. Ross, and Lu Jiang. Language model beats diffusion – tokenizer is key to visual generation, 2024.

Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers, 2022.

Supplementary

A Discrete regularizers

To select a simple but performant vector quantization method for our discrete latent space, we compare the classic vector quantization (VQ) van den Oord et al. [2018] without additional complexity like codebook reinitialization, lookup free quantization (LFQ) Yu et al. [2024], and finite scalar quantization (FSQ) Mentzer et al. [2023]. While training these autoencoders we observe interesting differences in training dynamics with multiple crossing points between FSQ and LFQ for certain metrics. We present those in Figure 8 where we can see that FSQ often takes the lead in the beginning phases of training but eventually gives it up to LFQ. We can also see that both of these methods outperform classic VQ which struggles without additional aids.

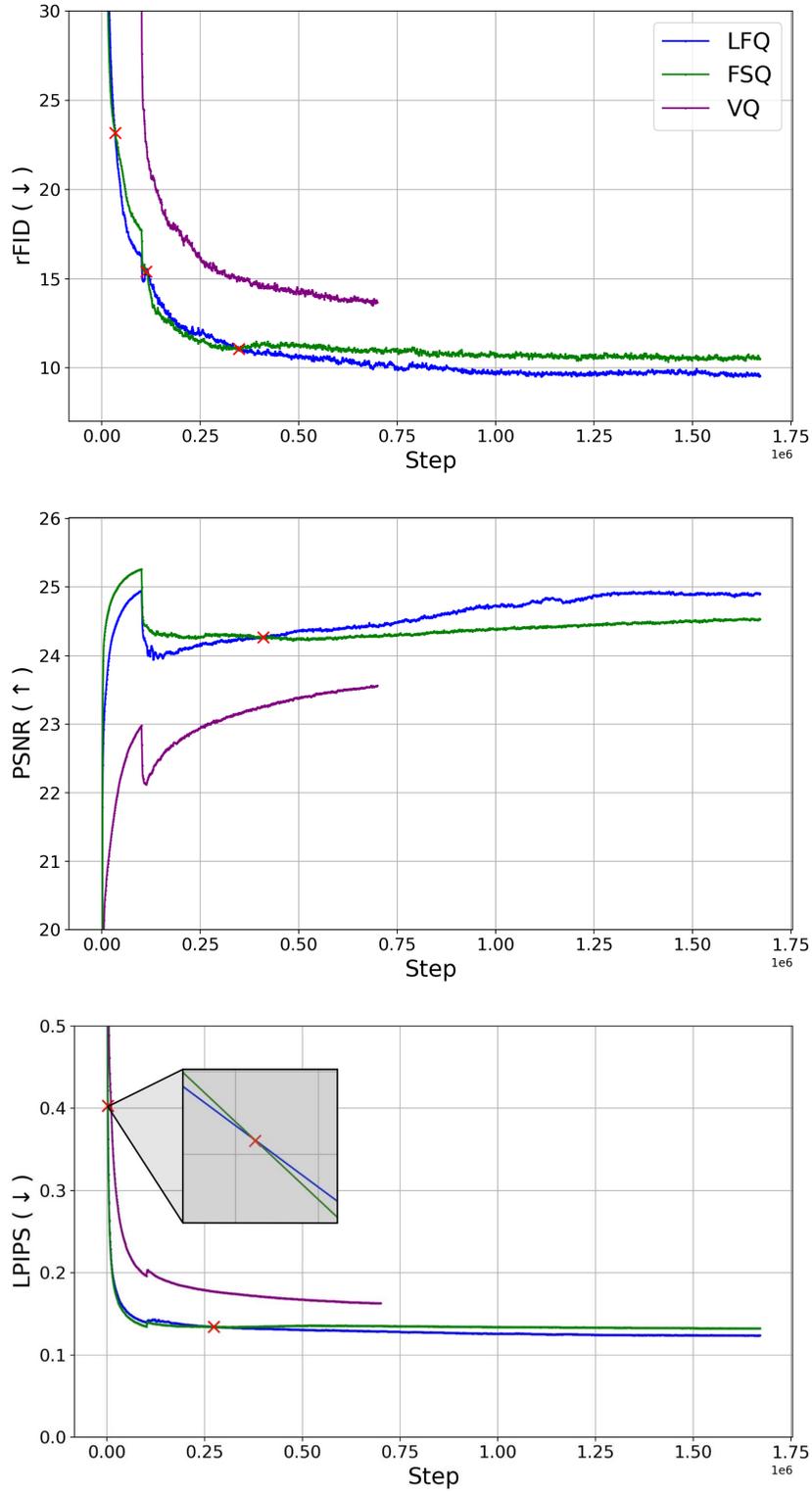


Figure 8: **Perceptual reconstruction metrics for various discrete regularization methods.** Classic vector quantization (VQ) struggles without tricks like codebook reinitialization. LFQ and FSQ have different training dynamics, often trading the lead in the beginning phases of training which is highlighted by the red X's.