# Causal Diffusion Transformers for Generative Modeling

Chaorui Deng    Deyao Zhu    Kunchang Li    Shi Guang    Haoqi Fan
ByteDance Research
https://github.com/causalfusion

## Abstract

*We introduce Causal Diffusion as the autoregressive (AR) counterpart of Diffusion models. It is a next-token(s) forecasting framework that is friendly to both discrete and continuous modalities and compatible with existing next-token prediction models like LLaMA and GPT. While recent works attempt to combine diffusion with AR models, we show that introducing sequential factorization to a diffusion model can substantially improve its performance and enables a smooth transition between AR and diffusion generation modes. Hence, we propose **CausalFusion** - a decoder-only transformer that dual-factorizes data across sequential tokens and diffusion noise levels, leading to state-of-the-art results on the ImageNet generation benchmark while also enjoying the AR advantage of generating an arbitrary number of tokens for in-context reasoning. We further demonstrate CausalFusion's multimodal capabilities through a joint image generation and captioning model, and showcase CausalFusion's ability for zero-shot in-context image manipulations. We hope that this work could provide the community with a fresh perspective on training multimodal models over discrete and continuous data.*

## 1. Introduction

Autoregressive (AR) and diffusion models are two powerful paradigms for data distribution modeling. AR models, also known as the next token prediction approach, dominate language modeling and are considered central to the success of large language models (LLMs) [5, 16, 46, 47, 61, 62]. On the other hand, diffusion models [13, 26, 29, 44], or score-based generative models [37, 54], have emerged as the leading approach for visual generation, driving unprecedented progress in the era of visual content generation [4, 17, 50].

The intrinsic distinction between AR and diffusion models lies in their approach to data distribution factorization. AR models treat data as an ordered sequence, factorizing it along the sequential axis, where the probability of each token is conditioned on all preceding tokens. This factorization enables the AR paradigm to generalize effectively and efficiently across arbitrary number of tokens, making it
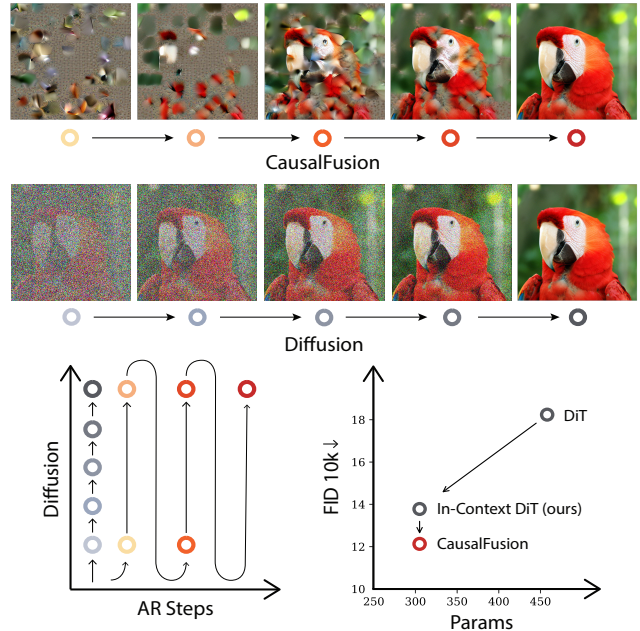


Figure 1. **Illustration of Dual-Factorization**. The arrow line indicates CausalFusion's generation path, moving from one state to the next by jointly generating along the sequential and noise-level dimension at each step. Compared to DiT, our In-context DiT substantially improves results with fewer parameters. CausalFusion further enhances performance without changing the architecture or parameter count. Results were trained on IN1K for 240 epochs. CausalFusion adopts arbitrary AR steps for image generation, but each step only diffuses partial tokens, resulting in similar (or slightly lower) computational complexity.

well-suited for long-sequence reasoning and in-context generation. In contrast, diffusion models factorize data along the noise-level axis, where the tokens at each step are a refined (denoised) version of themselves from the previous step. As a result, the diffusion paradigm is generalizable to arbitrary number of data refinement steps, enabling iterative quality improvement with scaled inference compute. While AR and diffusion models each excel within their respective domains, their distinct factorization approaches reveal complementary potential. Although recent studies [21, 72, 75] have attempted to integrate AR and diffusion within a sin-

1

(a) Samples generated by CausalFusion-XL/2, ImageNet 512×512, 800 epoch, DDPM 250 steps, CFG=4.0



(b) **Zero-shot image editing** results generated by CausalFusion-XL/2, ImageNet 512×512, 800 epoch. We first generate the original image (those on the left), then mask out its centre region, top-half, or bottom-half, and regenerate the image with new class conditions. Details are discussed in Sec 6.

Figure 2. **Visualization results**. All samples are generated by models trained only on **ImageNet-1K class-conditional generation** task, demonstrating CausalFusion's zero-shot image manipulation ability. See more visualization results in Appendix D.

gle model, they typically treat these paradigms as separate modes, missing the potential benefits of jointly exploring them within a 2-D factorization plane.

To this end, we introduce **CausalFusion**, a flexible framework that integrates both sequential and noise-level data factorization to unify their advantages. The degree of factorization along these two axes—namely, the AR step and diffusion step—is adjustable, enabling CausalFusion to revert seamlessly to the traditional AR or diffusion paradigms at either extreme. To enhance its generality, CausalFusion is designed to predict *any* number of tokens at *any* AR step, with *any* pre-defined sequence order and *any* level of inference compute, thereby minimizing the inductive biases presented in existing generative models. As shown in Figure 1, this approach provides a broad spectrum between the AR and diffusion paradigms, allowing smooth interpolation within two endpoints during both training and inference. Specifically, we explore CausalFusion in image generation and multimodal generation scenarios, where we observe that the level of training difficulties significantly influences the overall effectiveness of CausalFusion.

**Difficulties of generative tasks in CausalFusion:** Both AR and diffusion paradigms present unique challenges based on difficulties of their specific generative stages. In diffusion models, the effectiveness of training depends heavily on proper loss weighting across noise levels [22, 26], as higher noise levels are more difficult and usually

provide more valuable signals than lower noise levels. Similarly, AR models are susceptible to error accumulation [3] as early-stage predictions are made with limited visible context, making them more error-prone. Optimizing CausalFusion thus requires balancing across these varying task difficulties to optimize training signal impact and ensure sufficient exploration across the entire factorization plane.

In this paper, we formally examine the difficulties of generative tasks within CausalFusion. We show that, in addition to the noise levels in diffusion and the amount of visible context in AR, the total number of AR steps, which controls the interpolation between AR and diffusion, also plays a critical role in shaping training difficulties. Driven by these factors, we develop a scalable and versatile model based on the CausalFusion framework. Starting from the DiT architecture [44], we gradually convert it into a decoder-only transformer compatible with existing AR models like GPT [5, 46, 47] and LLaMA [16, 61, 62]. We provide insights on how to appropriately choose the number of AR steps during the training of CausalFusion models, and further introduce loss weighing along both the diffusion and AR axis to balance the impact of different generative stages. As shown in Figure 1 and 2, our model achieves state-of-the-art performance on the ImageNet class-conditional generation benchmark, significantly outperforming DiT [44] and enabling zero-shot image manipulations due to its AR nature. When pretraining on both text-to-image and image-

to-text tasks, our model surpasses forced-fusion frameworks such as TransFusion [75], demonstrating the versatility of our CausalFusion framework.

We highlight our main contribution below:

- We propose CausalFusion as the AR counterpart to DiT, achieving state-of-the-art results and enabling the unlimited token generation for in-context reasoning.
- We systematically study CausalFusion on the dual-factorization plane and identify key factors that improve the effectiveness of CausalFusion models.
- Compared with recent studies [75], CausalFusion enables a smooth, cohesive integration with language modeling for cross-modal generation and reasoning.

## 2. Related Works

**Diffusion Models**. Diffusion models [26, 52, 53] decompose the image generation task into a sequence of iterative denoising steps, gradually transforming noise into a coherent image. Early diffusion models [13, 43, 45, 49, 50] with U-net architectures pioneered denoising techniques for high-quality image synthesis. Later works [1, 44] like DiT shift from the U-net to transformer-based architectures, enabling greater compute scalability. Modern methods [8, 33] further extend DiT architectures leveraging significantly larger training resources to achieve impressive image generation quality.

**Autoregressive Generation**. Another popular approach for image generation involves autoregressive (AR) transformers that predict images token by token. Early works [14, 19, 48, 68] generated image tokens in raster order, progressing sequentially across the image grid. This rasterized approach was later identified as inefficient [6], prompting researchers to explore random-order generation methods [6, 7]. AR methods are further evolved to include new modalities such as video generation [31] and any-to-any generation [40, 70].

**Combining Diffusion and Autoregressive Models**. Recent models explore different methods for integrating AR and diffusion processes. DART [21] unifies AR and diffusion in a non-Markovian framework by conditioning on multiple historical denoising steps instead of only the current one. BiGR [23] generates discrete binary image codes autoregressively using a Bernoulli diffusion process. MAR [34] employs an AR model with a small diffusion head to enable continuous-value generation. Emu2 [57] applies an external diffusion module to decode its AR-based multimodal outputs. Compared to previous methods, CausalFusion focuses on autoregressive sequence factorization and decouples diffusion data processing across both sequential tokens and noise levels, achieving significant performance gains over traditional diffusion frameworks.

## 3. CausalFusion

**Preliminaries.** We first briefly review the Autoregressive (AR) and Diffusion paradigms in the context of image modeling before introducing our CausalFusion model. Both paradigms factorize the image distribution into a chain of conditional distributions. However, they do so along different axes. Given a sample of training images $\mathbf{X}$, AR models split $\mathbf{X}$ along the spatial dimensions into a sequence of tokens, $\mathbf{X} = \{\mathbf{x}_1, \ldots, \mathbf{x}_L\}$, where $L$ is the number of tokens. The joint distribution of $\mathbf{X}$ can be then factorized sequentially:

$$q(\mathbf{x}_{1:L}) = q(\mathbf{x}_1) \prod_{l=2}^{L} q(\mathbf{x}_l | \mathbf{x}_{1:l-1}). \tag{1}$$

During training, a neural network $p_\theta(\mathbf{x}_l | \mathbf{x}_{1:l-1})$ is trained to approximate $q(\mathbf{x}_l | \mathbf{x}_{1:l-1})$ by minimizing the cross-entropy $-\mathbb{E}_{q(\mathbf{x}_{1:L})} \log p_\theta(\mathbf{x}_{1:L})$. At inference time, the image is generated via the *next token prediction* paradigm.

In contrast, Diffusion models gradually add random noise (typically Gaussian) to $\mathbf{X}$ in a so-called *forward process*. It is a Markov chain along the noise level, where each noisy version $\mathbf{x}_t$ is conditioned on the previous state $\mathbf{x}_{t-1}$ as $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$. Here, $\beta_t$ is a variance schedule that ensures the forward process starts with a clean image $\mathbf{x}_0 = \mathbf{X}$ and gradually converges to random noise as $t \rightarrow T$. The joint distribution of $\mathbf{X}$ is then factorized as:

$$q(\mathbf{x}_{0:T}) = q(\mathbf{x}_0) \prod_{t=1}^{T} q(\mathbf{x}_t | \mathbf{x}_{t-1}). \tag{2}$$

To obtain $\mathbf{X}$ from noise, a neural network is trained to approximate the reverse transition of the forward process for $t \in [1, T]$:

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t), \Sigma_\theta(\mathbf{x}_t)) \tag{3}$$

As in AR models, training involves minimizing the cross-entropy between $q(\mathbf{x}_{0:T})$ and $p_\theta(\mathbf{x}_{0:T})$. In DDPM [26], $\Sigma_\theta(\mathbf{x}_t)$ is set to a constant value derived from the forward process, and $\mu_\theta(\mathbf{x}_t)$ is set to be the linear combination of $\mathbf{x}_t$ and a noise prediction model $\epsilon_\theta$ that predicts the noise $\epsilon$ of the forward process. This parameterization leads to the following training objective:

$$\min_\theta \mathbb{E}_{\mathbf{x}_0, \epsilon, t}[w(t) \| \epsilon - \epsilon_\theta(\mathbf{x}_t, t) \|^2] \tag{4}$$

where $w(t)$ is derived according to the noise schedule $\beta_t$, which gradually decays as $t$ grows. This objective is further simplified by setting $w(t) = 1$ for all $t$, results in a weighted evidence lower bound that emphasizes more difficult denoising tasks (i.e., larger noise level) at larger $t$ step.
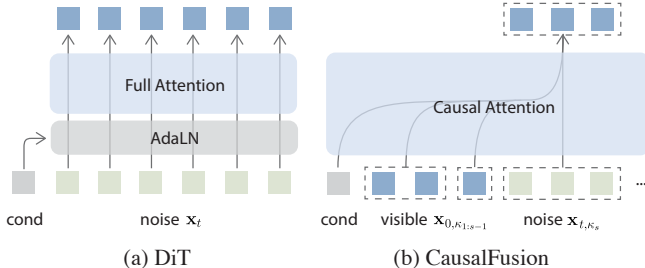
(a) DiT        (b) CausalFusion

Figure 3. **Conceptual comparison between the DiT and Causal-Fusion architectures**. a) DiT incorporates conditioning via adaptive layer normalization, processing a fixed-size set of entire image tokens as input. All the noise tokens $x_t$ are fed into DiT with full attention observation, enabling comprehensive modeling of the input during processing. b) CausalFusion treats all input modalities equally in an in-context manner, denoising a random subset of image tokens $x_{t,\kappa_s}$ at each step while causally conditioning on previously denoised tokens $x_{0,1:\kappa_s-1}$, and other contextual inputs. This approach enforces the model to reconstruct the image with partial observation, embodying the spirit of masked feature prediction models [24, 35, 67].

**Our approach.** From the above formulations, the AR and diffusion paradigms naturally support the scaling of sequence length and denoising steps, respectively, offering complementary advantages for image generation. To unify their advantages, we present CausalFusion, a general paradigm that scales effectively in both directions.

We start by directly extending Eqn. (2) to encompass the AR factorization:

$$q(\mathbf{x}_{0:T,\kappa_s}|\mathbf{x}_{0,\kappa_{1:s-1}}) =$$
$$q(\mathbf{x}_{0,\kappa_s}) \prod_{t=1}^{T} q(\mathbf{x}_{t,\kappa_s}|\mathbf{x}_{t-1,\kappa_s}, \mathbf{x}_{0,\kappa_{1:s-1}}) \quad (5)$$

for $s \in [1, S]$. Here, $S$ denotes the total number of AR steps, while $\kappa_s$ is an index set that identifies the subset of image tokens to be processed at the $s$-th AR step, with $|\kappa_s|$ representing the number of tokens in this subset. Each AR step processes only the tokens indicated by $\kappa_s$, isolating specific portions of the image, as shown in the top row of Figure 1. The term $\mathbf{x}_{t,\kappa_s}$ represents the dual-factorized image tokens at the $s$-th AR step and $t$-th diffusion step.

During training, the objective of our CausalFusion model is to approximate $p_\theta(\mathbf{x}_{t-1,\kappa_s}|\mathbf{x}_{t,\kappa_s}, \mathbf{x}_{0,\kappa_{1:s-1}})$ for all $t$ and $s$. Compared with the formulation in Eqn. (3), Causal-Fusion requires the training sequence to contain not only noised image tokens at the current AR step $\mathbf{x}_{t,\kappa_s}$, but also the clean image tokens from all previous AR steps $\mathbf{x}_{0,\kappa_{1:s-1}}$, allowing the model to leverage information from earlier AR steps to refine the current tokens effectively. A generalized version of causal attention mask is also required to prevent $\mathbf{x}_{0,\kappa_{1:s-1}}$ from observing $\mathbf{x}_{t,\kappa_s}$. During inference, the dual-factroization in Eqn. (5) enables CausalFusion to generate

| Model | Params (M) | FID10k↓ |
|---|---|---|
| <u>DiT</u> [44] | 458 | 18.24 |
| - AdaLN-zero [44] | 305 | 26.71 |
| + new recipe | 305 | 21.94 |
| + T embedding | 308 | 20.68 |
| + QK-norm | 308 | 18.66 |
| + lr warmup | 308 | 17.11 |
| + All (In-context DiT) | 308 | **13.78** |

Table 1. **In-context DiT baseline**. ImageNet 256×256, 240 epoch. Baseline settings are marked by <u>underlines</u> and selected settings highlighted in gray.

an unlimited sequence of image tokens through a *next token(s) diffusion* approach while enhancing the quality of each token by applying larger numbers of diffusion steps. See Figure 3(b) for an illustration of CausalFusion model architecture. Further details on the generalized causal attention mask can be found in Appendix A.

Notably, adhering to the principle of **minimal inductive bias**, CausalFusion imposes no restrictions on the number of AR steps $S$, the number of tokens processed at each AR step $|\kappa_s|$, or the specific token indices within each AR step. This flexibility enables a broad exploration space for generative modeling in both training and inference stages.

## 4. Initial studies on CausalFusion

To systematically study the design space of CausalFusion, we conduct experiments on the ImageNet dataset [12], where we train class-conditional image generation models at 256×256 resolution. We use the DiT-L/2 model as our base configuration. All models are trained with 240 epochs and evaluated using FID-10k (unless specified, we use FID-10K and FID interchangeably) and the ADM [13] codebase. Detailed training recipes and model configurations are provided in Appendix C.

**Baseline setup: In-context DiT.** As our target is to unify the AR and Diffusion paradigms, we need to unify their architectures first. To this end, we begin with the Transformer-based DiT [44] model. Following the DiT design, the 256×256 images are encoded into 32×32 latent representations [50] using a pretrained VAE model, followed by a 2×2 patchify layer that produces a sequence of $L = 256$ latent tokens. In original DiT, conditional information (e.g., label class) and the diffusion time step are incorporated through AdaLN-zero components, which are incompatible with decoder-only LLMs. To address this limitation, we adopt the in-context design of DiT from [44], treating the class and time step conditions as tokens and directly append them to the image token sequence. By default, we use 4 class tokens and one time step token. As a byproduct, this modification reduces the model size of in-context DiT to approximately $\frac{2}{3}$ of the AdaLN-zero version.
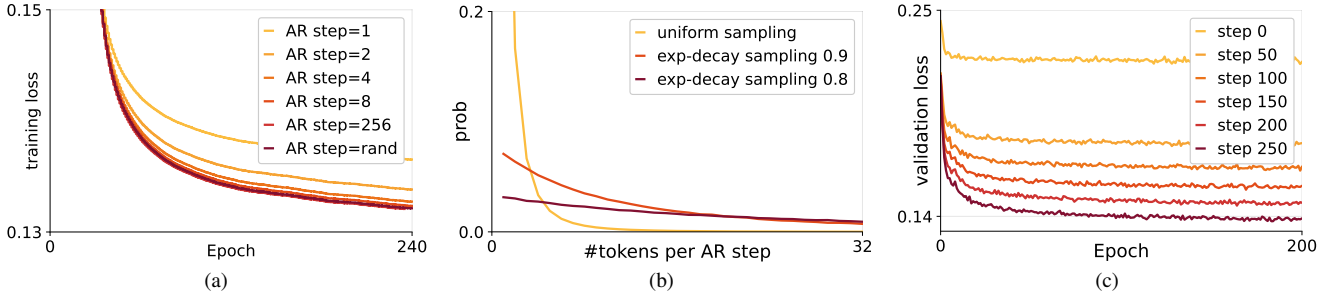
4

Figure 4. (a) Training loss using different number of AR steps. (b) Distribution of $|\kappa_s|$. (c) Validation loss at difference AR steps.

| #AR steps | FID10k↓ | | | |
| --- | --- | --- | --- | --- |
| | $S_{\text{eval}} = 1$ | $S_{\text{eval}} = 2$ | $S_{\text{eval}} = 4$ | $S_{\text{eval}} = 8$ |
| $S_{\text{train}} = 1$ | **13.78** | 356.69 | 404.67 | 390.18 |
| $S_{\text{train}} = 2$ | 16.69 | **14.77** | 47.49 | 136.04 |
| $S_{\text{train}} = 4$ | 24.14 | 15.37 | **18.13** | 33.14 |
| $S_{\text{train}} = 8$ | 54.08 | 24.49 | 22.66 | **20.01** |
| $S_{\text{train}} = 256$ | 313.28 | 321.62 | 261.26 | 192.25 |
| random | 21.31 | 22.17 | 23.54 | 25.05 |

Table 2. **Ablations on AR steps**. $S_{\text{train}}$ and $S_{\text{eval}}$ indicates the fixed AR steps used during training and inference, respectively. Baseline settings are marked by underlines and selected settings highlighted in gray.

To accelerate training, we use large batch sizes (e.g., 2048) and implement several improvements to stabilize training: (1) injecting the diffusion time step by adding a time step embedding to the image token embeddings instead of using a time step token; (2) applying head-wise QK layer normalization within the self-attention layers, following the practices in [11]; and (3) incorporating a learning rate warmup stage during training.

The impact of our new designs is shown in Table 1. Initially, the native In-context DiT from [44] performs significantly worse than the AdaLN-zero version. Our revised training recipe improves performance to 21.94 FID. Incorporating the time step embedding and head-wise QK norm further enhances performance, achieving an FID of 18.66. Adding the learning rate warmup yields an additional improvement. Overall, our final in-context DiT-L/2 model, though conceptually simple, reaches an FID-10k of 13.78—competitive to the best-performing DiT-XL/2 model (12.92 FID-10k) in [44]—and serves as a robust baseline that can be steadily trained with large batch sizes.

**CausalFusion with fixed number of AR steps.** Building on the In-context DiT baseline, we begin with a simplified version of CausalFusion that uses a fixed number of AR steps $S$ during both training and inference, with the number of tokens predicted at each AR step fixed to $|\kappa_s| = \frac{L}{S}$. Specifically, we modify the input sequence to include clean image tokens and use generalized causal attention masks within the attention modules. Figure 3 illustrates the architectural differences between DiT and CausalFusion. We

train several CausalFusion models with different numbers of AR steps, i.e., $S = 1, 2, 4, 8$, and 256. Here, $S = 1$ indicates the In-context DiT, while $S = 256$, equivalent to $L$, represents the pure AR training mode. To evaluate these models, we first use the same number of AR steps as in training, and further study their generalization performance to other number of AR steps.

As shown in Table 2, CausalFusion trained with fixed AR steps cannot be robustly transferred to other inference settings. E.g., all models yield substantially worse performance when their inference settings are not aligned with training. By comparing the best evaluation result of each training setting, we observe that increasing the number of AR steps leads to a huge decline in performance. Specifically, the 8-step CausalFusion yields an FID of 20.01, clearly lagging behind the 13.78 FID achieved by In-context DiT. However, from the loss curves in Figure 4(a), an opposite trend is observed, where models trained with more AR steps consistently exhibit lower loss values than those with fewer AR steps. This suggests that the learning tasks become over-simplified as the number of AR steps increases.

**CausalFusion with random number of AR steps.** Additionally, we train a CausalFusion model where the number of AR steps $S$ are uniformly sampled from 1 to $L$, with $|\kappa_s|$ also randomly set at each AR step. We evaluate this model across various inference settings, same as above. As shown in Table 2, this training setting performs relatively consistent under different inference AR steps compared to those trained with fixed AR steps, demonstrating its greater flexibility during training and versatility during inference. However, this setting still leads to inferior results compared to the In-context DiT baseline, e.g., an FID of 21.31 when evaluated with a single AR step ($S = 1$) as diffusion mode. Figure 4(b) offers further insight into this behavior, showing that uniform AR step sampling during training leads to a highly imbalanced $|\kappa_s|$ distribution. As a result, the training signal is dominated by AR steps with very few tokens—over 95% of AR steps have $|\kappa_s| \leq 16$. These steps are uniformly distributed along the AR axis, causing the model to overly rely on visible context and thus diminishing the complexity of the training task.

5

| ratio | FID10k↓ | | | |
|---|---|---|---|---|
| | $S_{eval}=1$ | $S_{eval}=2$ | $S_{eval}=4$ | $S_{eval}=8$ |
| <u>1.0</u> | 21.31 | 22.17 | 23.54 | 25.05 |
| 0.95 | 14.49 | 17.78 | 19.79 | 23.93 |
| 0.9 | 12.89 | 15.57 | **18.83** | **22.72** |
| 0.85 | 12.94 | 15.54 | 19.12 | 23.46 |
| 0.8 | **12.78** | **15.42** | 19.38 | 23.78 |

| Patch order | FID10k↓ |
|---|---|
| raster order | 14.46 |
| block-wise raster (8x8) | 14.76 |
| block-wise raster (4x4) | 14.62 |
| dilated order | 15.54 |
| random order | **12.89** |

| weight | FID10k↓ | | |
|---|---|---|---|
| | $S_{eval}=1$ | $S_{eval}=2$ | $S_{eval}=4$ |
| <u>1→1</u> | 12.89 | 15.57 | 18.83 |
| 1.5→1 | 12.61 | 15.49 | 18.32 |
| 2→1 | **12.13** | **15.15** | 18.09 |
| 2.5→1 | 12.32 | 15.22 | 17.99 |
| 3→1 | 12.50 | 15.28 | **17.92** |

(a) **Exponential decay** in AR step sampling. A proper decay ratio leads to competitive or better performance across all inference settings than using fixed AR steps.

(b) **Token order** influences the locality of image tokens and further affects training difficulty.

(c) **AR loss weighting** boosts performance by facilitating better learning from difficult samples.

Table 3. **Ablations on AR step decay, ordering, and AR weighting**. Baseline settings are marked by <u>underlines</u> and selected settings highlighted in  gray .

Lastly, given the CausalFusion model trained with random AR steps, we compare the loss values produced by different AR steps on the validation set. As shown in Figure 4(c), later AR steps yield much lower loss values than earlier steps, suggesting a clear trend of vanished training signals along the AR axis.

## 5. Shaping task difficulties in CausalFusion

Based on the above observations, we aim to adjust the difficulties of the generative tasks in CausalFusion to balance training signal impact and ensure thorough exploration of the factorization space. By default, we use random AR steps during training due to its effectiveness in generalizing across various inference settings. Building on this setup, we identify several straightforward approaches that effectively shape task difficulties within CausalFusion, leading to significant performance improvements. We categorize the discussion into two parts: design choices for AR step sampling and loss weighting along the AR axis.

**Random AR steps with decayed sampling.** Instead of uniformly sampling the number of AR steps $S$ from $[1, L]$, we propose to exponentially decrease the sampling probability as $S$ increases, which alleviates the problem of imbalanced $|\kappa|_s$ distribution, as shown in Figure 4(b). As a result, large $|\kappa_s|$ appears more frequently in the training sequence, and more tokens are predicted base on less visible context. We introduce a hyper-parameter $\gamma \leq 1$ to control the exponential decay ratio where $\gamma = 1$ denotes the naive CausalFusion model trained with uniformly sampled AR steps, while $\gamma = 0$ represents our In-context DiT baseline. From Table 3(a), by decreasing $\gamma$ from 1.0 to 0.95, we obtain remarkable performance improvements across all inference settings, with gains of nearly 7 and 5 points when $S_{eval} = 1$ and 2, respectively. Furthermore, when $\gamma = 0.9$, CausalFusion surpasses the strong In-context DiT using the pure diffusion inference mode, and performance in other inference settings is further enhanced. While smaller values of $\gamma$, such as 0.8, yield even better performance with one AR step evaluation, we set the defualt value to 0.9 as it provides a balanced improvement across all inference settings.

**Loss weighting along AR axis.** We modify the weighting term $w(\cdot)$ in Eqn. (4) to further consider the AR step $s$. In practice, we simply set $w(s, t)$ to a pre-defined value $\lambda \geq 1$ at $s = 1$, and linearly decay it to 1 at $s = S$, and keep using constant weight for different $t$. In this way, the model is trained to focus more on the hard generative tasks at early AR steps or larger noise levels. We analysis the impact of $\lambda$ in Table 3(c). From the table, setting $\lambda$ to a proper value improves the performance. Intuitively, as the model generates closer to the end of the AR axis, the task becomes easier due to high **locality**[65] in the visible context, causing some generative tasks to degrade into local feature interpolation. In contrast, predictions made during earlier AR steps facilitate the learning of non-local dependencies within the visual context, which is beneficial to generative modeling.

**Difficulty vs. locality.** The hypothesis of locality aligns with our observations in Table 3(b), where using random sequence order in CausalFusion during training significantly outperforms manually assigned orders. Specifically, using fixed (block-wise) raster order leads the model to rely excessively on local tokens, which makes the training task easier. In contrast, CausalFusion is trained with a random order by default, following the principle of *minimal inductive bias*. This design encourages the model to develop robust generative modeling abilities rather than relying on fixed ordering priors, while allowing flexible inference orders.

## 6. Performance comparison

**Class-conditional image generation.** We evaluate our final method on the ImageNet class-conditional generation benchmark. For system-level comparisons, we use 64 tokens to encode the class condition. The impact of varying the number of class tokens is analyzed in Appendix B. We train three sizes of CausalFusion models: CausalFusion-L (368M), CausalFusion-XL (676M), and CausalFusion-H (1B). All models are trained for 800 epochs with a batch size of 2048. By default, we use a single AR step inference with 250 DDPM steps as in DiT [44], and report FID-50k for benchmarking against existing models. Detailed hyperparameters are provided

| | | 256×256, w/o CFG | | | | 256×256, w/ CFG | | | | 512×512, w/ CFG | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Params | FID↓ | IS↑ | Pre.↑ | Rec.↑ | FID↓ | IS↑ | Pre.↑ | Rec.↑ | FID↓ | IS↑ | Pre.↑ | Rec.↑ |
| GIVT [63] | 304M | 5.67 | - | 0.75 | 0.59 | 3.35 | - | 0.84 | 0.53 | 2.92 | - | 0.84 | 0.55 |
| MAR-B [34] | 208M | 3.48 | 192.4 | 0.78 | 0.58 | 2.31 | 281.7 | 0.82 | 0.57 | - | - | - | - |
| LDM-4 [50] | 400M | 10.56 | 103.5 | 0.71 | 0.62 | 3.6 | 247.7 | **0.87** | 0.48 | - | - | - | - |
| CausalFusion-L | 368M | **5.12** | **166.1** | **0.73** | **0.66** | **1.94** | **264.4** | 0.82 | **0.59** | - | - | - | - |
| MAR-L [34] | 479M | 2.6 | 221.4 | 0.79 | 0.60 | 1.78 | 296.0 | 0.81 | 0.60 | 1.73 | 279.9 | - | - |
| ADM [13] | 554M | 10.94 | - | 0.69 | 0.63 | 4.59 | 186.7 | 0.82 | 0.52 | 3.85 | 221.7 | 0.84 | 0.53 |
| DiT-XL [44] | 675M | 9.62 | 121.5 | 0.67 | 0.67 | 2.27 | 278.2 | **0.83** | 0.57 | 3.04 | 240.8 | 0.84 | 0.54 |
| SiT-XL [42] | 675M | 8.3 | - | - | - | 2.06 | 270.3 | 0.82 | 0.59 | 2.62 | 252.2 | **0.84** | 0.57 |
| ViT-XL [22] | 451M | 8.10 | - | - | - | 2.06 | - | - | - | - | - | - | - |
| U-ViT-H/2 [1] | 501M | 6.58 | - | - | - | 2.29 | 263.9 | 0.82 | 0.57 | 4.05 | - | - | - |
| MaskDiT [73] | 675M | 5.69 | 178.0 | 0.74 | 0.60 | 2.28 | 276.6 | 0.80 | 0.61 | 2.50 | 256.3 | 0.83 | 0.56 |
| RDM [59] | 553M | 5.27 | 153.4 | 0.75 | 0.62 | 1.99 | 260.4 | 0.81 | 0.58 | - | - | - | - |
| CausalFusion-XL | 676M | **3.61** | **180.9** | **0.75** | **0.66** | **1.77** | **282.3** | 0.82 | **0.61** | **1.98** | **283.2** | 0.83 | **0.58** |

Table 4. **System performance comparison** on ImageNet class-conditioned generation. Numbers marked with gray blocks use temperature sampling during inference.

| | Type | Tokenizer | Params | Training Epoch | Sampler (Steps) | Sampling tricks | FID↓ |
|---|---|---|---|---|---|---|---|
| Open-MAGVIT2-L [41] | AR | MAGVIT2 | 800M | 300 | AR(256) | N/A | 2.51 |
| Open-MAGVIT2-XL [41] | AR | MAGVIT2 | 1.5B | 300 | AR(256) | N/A | 2.33 |
| LlamaGen-3B [56] | AR | custom | 3.1B | - | AR(256) | N/A | 2.18 |
| VAR-d24 [60] | VAR | custom | 1B | 350 | VAR | N/A | 2.09 |
| VAR-d30 [60] | VAR | custom | 2B | 350 | VAR | reject sampling | 1.73 |
| Simple-diffusion [27] | Diffusion | N/A | 2B | 800 | DDPM | N/A | 2.44 |
| FiTv2-3B [66] | Diffusion | SD | 3B | 256 | DDPM(250) | N/A | 2.15 |
| VDM++ [30] | Diffusion | N/A | 2B | - | EDM | - | 2.12 |
| Large-DiT-7B [20] | Diffusion | SD | 3B | 435 | DDPM(250) | N/A | 2.10 |
| Flag-DiT-3B [20] | Diffusion | SD | 3B | 256 | adaptive Dopri-5 | N/A | 1.96 |
| DiT-MoE-XL/2-8E2A [18] | Diffusion | SD | 16B | ≈1000 | DDPM(250) | N/A | 1.72 |
| DiMR-G/2R [38] | Diffusion | SD | 1.1B | 800 | DPM-solver(250) | N/A | 1.63 |
| DART-XL [21] | AR+Diffusion | LDM | 812M | - | AR(256)+FM(100) | $\tau$ sampling | 3.98 |
| MonoFormer [72] | AR+Diffusion | SD | 1.1B | - | DDPM(250) | N/A | 2.57 |
| BiGR-XL-d24 [23] | AR+Diffusion | custom | 799M | 400 | AR(256)+DDPM(100) | $\tau$ sampling | 2.49 |
| BiGR-XXL-d32 [23] | AR+Diffusion | custom | 1.5B | 400 | AR(256)+DDPM(100) | $\tau$ sampling | 2.36 |
| MAR-H [34] | AR+Diffusion | custom | 943M | 800 | AR(256)+DDPM(100) | $\tau$ sampling | 1.55 |
| CausalFusion-H | Diffusion | custom | 1B | 800 | DDPM(250) | N/A | 1.64 |
| CausalFusion-H | Diffusion | custom | 1B | 800 | DDPM(250) | CFG interval | 1.57 |

Table 5. **System performance comparison** on 256×256 ImageNet generation, compared with previously reported large models.

in Appendix C. As shown in Table 4, on 256×256 image generation, CausalFusion-L achieves an FID-50k of 5.12 without classifier-free guidance [25] (CFG), outperforming DiT-XL/2 by 4.5 points with 50% fewer parameters. CausalFusion-XL further improves this result to 3.61, and when using CFG, achieves a state-of-the-art result of 1.77, significantly outperforming strong baselines like DiT and SiT. Additionally, CausalFusion-XL demonstrates effectiveness in high-resolution generation, achieving an FID of 1.98 on 512×512 images with CFG.

We also provide a system-level comparison with existing methods in Table 5. CausalFusion-H achieves an FID of 1.64 using the standard 250-step DDPM sampler, outperforming previous diffusion models with larger model sizes, such as FiTv2-3B [66] and Large-DiT-7B [20], and achieving comparable results to DiMR-G/2R (1.64 vs. 1.63) despite DiMR-G/2R's use of a stronger sampler (DPM-solver [51]). By applying the CFG interval [32] approach, CausalFusion-H further improves to an FID of 1.57, positioning it among the top-performing models on the ImageNet 256×256 benchmark.



A bear wearing a chef's hat baking cookies in a log cabin.

A cactus wearing a sombrero and sunglasses in the desert.

A goose in rain boots, stomping through puddles in the park.

A cupcake with sprinkles, lounging under a tiny umbrella on a beach.

(a) Samples on Text-to-Image generation.



A man in a red jacket riding a bicycle down a city street.

A woman in a yellow shirt is standing next to a large elephant at night.

A kitchen with a stainless steel dishwasher and a wooden cabinet.

A cat sitting on a red and white chair.

(b) Samples on Image Caption generation.

Figure 5. **Multimodal generation**. Results are generated by a single CausalFusion XL model trained on ImageNet recaption data.

|  | Source | Size | FID30k↓ | CIDEr↑ |
|---|---|---|---|---|
| Transfusion-L [75] | IN1KCap | 1M | 8.1 | 34.5 |
| CausalFusion-L | IN1KCap | 1M | **7.1** | **47.9** |

(a) MSCOCO [36] is used for FID30k and CIDEr evalution.

|  | Params | Data | Size | FID10k↓ | Acc↑ | CIDEr↑ |
|---|---|---|---|---|---|---|
| DiT [44] | 458M | IN1K | 1M | 18.2 | 83.5 | 94.4 |
| CausalFusion | 368M | IN1K | 1M | 11.8 | 84.2 | 98.0 |
| CausalFusion† | 368M | IN1K | 1M | **9.3** | **84.7** | **103.2** |

(b) ImageNet is used for FID10k and Acc, MSCOCO is used for CIDEr evalution.

Table 6. (a) **Comparison with Transfusion** [75] on perception and generation benchmarks. All models are trained under the same settings using the same pretraining data. (b) **Comparison with DiT** [44] on perception and generation benchmarks. The model marked with † is trained with a VAE from [34], using a loss function to predict latent variables rather than noise.

**Zero-shot image editing.** CausalFusion naturally supports zero-shot image editing, as it is trained to predict a random subset of image tokens based on a random subset of visible image tokens. This inherent flexibility enables the model to perform localized edits without requiring task-specific fine-tuning. As shown in Figure 2(b), our model can generate high-quality editing results even when only pretrained on the ImageNet class-conditional generation task, demonstrating its robustness and adaptability to editing tasks. Moreover, CausalFusion's dual-factorized design allows it to balance contextual coherence with high-fidelity updates, ensuring that edited regions blend seamlessly into the surrounding content. See Appendix D for additional visualizations showcasing the model's ability to handle diverse editing scenarios.

**Vision-Language joint modeling.** CausalFusion can integrate the language modality by applying a separate next-token prediction loss on text, similar to GPT [46], enabling it to jointly model both image and text data. In this experiment, CausalFusion was trained on two tasks simultaneously: Text-to-Image (T2I) generation and Image Captioning. During training, text precedes the image in 90% of cases, framing it as a T2I task where only the image loss is applied. In the remaining cases, the text follows the image, and both text loss (for Image Captioning) and image loss (for classifier-free guidance [25] in T2I) are applied, with the text loss weighted at 0.01 relative to the image loss.

We follow the configurations and training/inference protocols from previous sections. For language tokenization, we use the LLaMA-3 [16] tokenizer. Models are trained on a re-captioned ImageNet dataset, with each image labeled by 10 captions generated by Qwen2VL-7B [64]. T2I and Image Captioning tasks are evaluated using zero-shot MSCOCO-30k FID and zero-shot MSCOCO CIDEr on Karpathy's test split, respectively. We compare Causal-Fusion with a contemporary multimodal model, Transfusion [75], which integrates language and vision modeling using standard diffusion loss for images and next-token prediction loss for text. In Transfusion, language tokens are conditioned on image embeddings with added diffusion noise. As Transfusion is not open-sourced, we implemented it based on the original paper, aligning the model architecture, VAE encoder, and language tokenizer with those used in CausalFusion. Results with 240-epoch training are

presented in Table 6(a). Compared to Transfusion, Causal-Fusion demonstrates superior performance in both text-to-image generation and image captioning, highlighting its strong potential as a foundational model for multimodal tasks. In Figure 5, we show a single pretrained CausalFu-sion XL model performing text-to-image generation at the top and image-to-text generation (image captioning) at the bottom. Further experiment details, including data, model design, and hyperparameters, are provided in Appendix C.

**Visual Representation Learning.** We further evaluate CausalFusion models from a representation learning perspective. Specifically, we leverage the CausalFusion model pretrained on the 256×256 ImageNet class-conditional generation task and fine-tune it on the ImageNet classification and MSCOCO captioning tasks. For image classification, we use the average-pooled features from the last layer of CausalFusion, followed by a linear classifier. For image captioning, we add a small transformer encoder-decoder module as the language head on top of CausalFusion. The pretrained CausalFusion model follows the default configuration described in previous sections. We compare it to the DiT-L/2 model pretrained for the same number of epochs. Detailed hyperparameters for fine-tuning are provided in Appendix C. As shown in Table 6(b), our CausalFusion model outperforms DiT on all fine-tuning tasks, indicating that CausalFusion learns superior representations compared to DiT. We hypothesize that the random grouped token diffusion mechanism in CausalFusion, which diffuses images with partially observed inputs, implicitly performs as masked representation prediction [24, 67], enhancing the model's representation learning ability.

## 7. Conclusion

We propose CausalFusion, a decoder-only transformer that unifies AR and diffusion paradigms through a dual-factorized framework across sequential tokens and diffusion noise levels. This approach achieves state-of-the-art performance on the ImageNet generation benchmark, supports arbitrary-length token generation, and enables smooth transitions between AR and diffusion modes. CausalFusion also demonstrates multimodal capabilities, including joint image generation and captioning, as well as zero-shot image manipulations. Our framework offers a new perspective on unified learning of diffusion and AR models.

## A. Generalized Causal Attention

We design the Generalized Causal Attention for CausalFusion models. The core idea is to maintain causal dependencies across all AR steps while ensuring that each AR step relies only on *clean* image tokens from preceding AR steps. This design allows CausalFusion to generate images using a *next-token(s) diffusion* paradigm. We show an example of the generalized causal attention mask in Figure 6, and the PyTorch-style pseudo code for obtaining generalized causal mask in Algorithm 1.



Figure 6. **Generalized causal mask.** In this case, the input sequence is organized to have 3 AR-steps $\kappa_1$, $\kappa_2$, and $\kappa_3$, containing 2, 2, and 3 tokens, respectively. $\mathbf{x}_{0,\kappa_1}$ and $\mathbf{x}_{0,\kappa_2}$ are the clean tokens at the first two AR steps, while $\mathbf{x}_{t,\kappa_1}$, $\mathbf{x}_{t,\kappa_2}$, and $\mathbf{x}_{t,\kappa_3}$ are noised tokens. White and gray blocks denote the masked and unmasked attention, respectively. Note that, each $\mathbf{x}_{t,\kappa_s}$ only attends to itself and the clean tokens from previous AR steps $\mathbf{x}_{0,\kappa_{1:s-1}}$.

## B. More Analyses

**Diffusion time steps sampling.** Following DiT [44], we randomly sample the diffusion time step $t$ during training. By default, the same $t$ is used across all AR steps when training CausalFusion models. Here, we explore the impact of using different $t$ values for different AR steps during training. The training and evaluation settings remain consistent with Section 4. As shown in Table 7, using either shared or random $t$ values results in similar performance, indicating that CausalFusion is robust to this variation. Additionally, we evaluate a setting where multiple diffusion time steps are sampled for each AR step. Specifically, we experiment with sampling 4 and 8 different time steps, reducing the batch size by factors of $4\times$ and $8\times$, respectively, to ensure the total number of tokens used for loss calculation remains constant. As shown in the table, using multiple diffusion time steps per AR step achieves comparable performance to the default setting, further demonstrating the robustness of CausalFusion to this factor. Notably, in this approach, the clean image tokens $\mathbf{x}_{0,\kappa_s}$ at each AR step need

**Algorithm 1** Generalized causal mask

```
def get_attn_mask(ctx_len, x_len, step):
    # tx_len: the length of clean tokens
    # x_len: the length of noisy tokens
    # step: number of AR steps

    # sample random tokens per AR step
    sumk = random.sample(range(1, x_len), step - 1)
    sumk = [0] + sorted(sumk) + [x_len]

    # build `causal` masks
    seq_len = ctx_len + x_len
    attn_mask = torch.ones(size=(seq_len, seq_len))
    m1 = torch.ones(size=(ctx_len, ctx_len))
    m2 = torch.ones(size=(x_len, ctx_len))
    m3 = torch.ones(size=(x_len, x_len))
    for i in range(len(sumk) - 2):
        m1[sumk[i]:sumk[i+1], 0:sumk[i+1]] = 0
        m2[sumk[i+1]:sumk[i+2], 0:sumk[i+1]] = 0
    for i in range(len(sumk) - 1):
        m3[sumk[i]:sumk[i+1], sumk[i]:sumk[i+1]] = 0

    attn_mask[:ctx_len, :ctx_len] = m1
    attn_mask[ctx_len:, :ctx_len] = m2
    attn_mask[ctx_len:, ctx_len:] = m3
    return attn_mask # 1 for mask, 0 for unmask
```

to be computed only once and can be shared across multiple $\mathbf{x}_{t,\kappa_s}$ with different $t$ values. Consequently, the additional computational cost introduced by clean image tokens during training is minimized, amounting to only $\sim$10%.

|  | FID10k |
|---|---|
| shared $t$ for different AR steps | 12.13 |
| random $t$ for different AR steps | 12.27 |
| $4\times$ $t$ for each AR step | 12.19 |
| $8\times$ $t$ for each AR step | 12.23 |

Table 7. **Diffusion time steps sampling** strategy does not affect the performance. The default setting is underlined.

| #class tokens | params (M) | FID10k |
|---|---|---|
| 4 | 308 (+3.9) | 12.13 |
| 16 | 320 (+15.6) | 12.04 |
| 64 | 368 (+62.5) | 11.84 |
| 1 (repeat 64$\times$) | 305 (+1.0) | 12.29 |
| 4 (repeat 16$\times$) | 308 (+ 3.9) | **11.75** |

Table 8. **#Class tokens** offers a trade-off between performance and number of parameters. The default setting is underlined.

**Class condition tokens.** As discussed in Sections 4 and 6, we use 4 class condition tokens for ablation studies and 64 tokens for system-level comparisons. Here, we examine the impact of the number of class tokens in the CausalFusion framework. As shown in Table 8, increasing the number of class tokens to 64 slightly improves performance (12.13 vs. 11.84 FID). However, this also adds 62.5M parameters, a significant increase (20%) for a CausalFusion-L model with 304M parameters. To address this, we adopt a token-repeating strategy from [34], which achieves comparable performance (11.75 vs. 11.84 FID) without increasing the parameter count. This finding suggests that the computation allocated to class conditioning is more critical than the number of parameters dedicated to it.

# C. Implementation Details

**Class-conditional image generation.**    In Table 9, we provide the detailed settings of CausalFusion models for class-conditional image generation in Section 4 and 5.

| config | value |
|---|---|
| image resolution | 256×256 |
| hidden dimension | 1024 |
| #heads | 16 |
| #layers | 24 |
| #cls tokens | 4 |
| patch size | 2 |
| positional embedding | sinusoidal |
| VAE | SD [55] |
| VAE donwsample | 8× |
| latent channel | 4 |
| optimizer | AdamW [39] |
| base learning rate | 1e-4 |
| weight decay | 0.0 |
| optimizer momentum | $\beta_1, \beta_2 = 0.9, 0.95$ |
| batch size | 2048 |
| learning rate schedule | constant |
| warmup epochs | 40 |
| training epochs | 240 |
| augmentation | horizontal flip, center crop |
| diffusion sampler | DDPM [26] |
| diffusion steps | 250 |
| evaluation suite | ADM [13] |
| evaluation metric | FID-10k |

Table 9. Ablation study configuration.

**System-level comparisons.**    In Table 10, we provide the detailed settings of CausalFusion models for system-level comparisons in Section 6.

| config | value |
|---|---|
| hidden dimension | 1024 (L), 1280 (XL), 1408 (H) |
| #heads | 16 (L), 20 (XL), 22 (H) |
| #layers | 24 (L), 32 (XL), 40 (H) |
| #cls tokens | 64 |
| positional embedding | learnable |
| VAE | mar [34] |
| VAE donwsample | 16× |
| latent channel | 16 |
| optimizer | AdamW [39] |
| base learning rate | 1e-4 |
| weight decay | 0.0 |
| optimizer momentum | $\beta_1, \beta_2 = 0.9, 0.95$ |
| batch size | 2048 |
| learning rate schedule | constant |
| warmup epochs | 40 |
| training epochs | 800 |
| augmentation | horizontal flip, center crop |
| diffusion sampler | DDPM [26] |
| diffusion steps | 250 |
| evaluation suite | ADM [13] |
| evaluation metric | FID-50k |

Table 10. System-level comparison configuration.

**Multi-modal CausalFusion.**    In Table 11, we provide the detail experiment hyperparameters for both CausalFusion and Transfusion experiments in Section 6. The training

dataset is augmented with 10 captions per image from ImageNet, generated using Qwen2VL-7B-Instruct [64] with the following prompt:

*You are an image captioner. You need to describe images in COCO style. COCO style is short. Here are some examples of COCO style descriptions: 'A car that seems to be parked illegally behind a legally parked car' 'This is a blue and white bathroom with a wall sink and a lifesaver on the wall.' 'Meat with vegetable and fruit displayed in roasting pan.' 'Group of men playing baseball on an open makeshift field.'*

| config | value |
|---|---|
| image resolution | 256×256 |
| hidden dimension | 1024 |
| #heads | 16 |
| #layers | 24 |
| #max text tokens | 35 |
| patch size | 2 |
| image positional embedding | sinusoidal |
| text positional embedding | learnable |
| VAE | SD [55] |
| VAE donwsample | 8× |
| latent channel | 4 |
| optimizer | AdamW [39] |
| base learning rate | 1e-4 |
| text loss coefficient | 0.01 |
| weight decay | 0.0 |
| optimizer momentum | $\beta_1, \beta_2 = 0.9, 0.95$ |
| batch size | 2048 |
| learning rate schedule | constant |
| warmup epochs | 40 |
| training epochs | 240 |
| augmentation | horizontal flip, center crop |
| diffusion sampler | DDPM [26] |
| diffusion steps | 250 |
| generation eval. metric | MSCOCO 0-shot FID-30k |
| captioning eval. metric | MSCOCO CIDEr (Karpathy test) |

Table 11. Multi-modal experiment configuration for both Causal-Fusion and Transfusion.

**Fine-tuning for ImageNet classification.**    When fine-tuning our CausalFusion model for ImageNet classification, we adhere to the basic architecture of the Vision Transformer (ViT) [15], with the exception of the class token. We exclude the extra timestep embedding, label embedding, and conditional position embedding. Layer normalization and a linear classification layer are applied to the averaged output tokens. Regarding hyperparameters, we follow the MAE training recipe [24] as detailed in Table 12, but we use BFloat16 precision during training to enhance stability.

**Fine-tuning for MSCOCO captioning.**    We follow the COCO caption fine-tuning setup of FLIP [35], incorporating an additional caption head consisting of a 3-layer transformer encoder and a 3-layer transformer decoder (with a width of 384 and 6 attention heads). This caption head takes

| config | value |
| --- | --- |
| optimizer | AdamW |
| base learning rate | 1e-3 (L) |
| weight decay | 0.05 |
| optimizer momentum | $\beta_1, \beta_2{=}0.9, 0.999$ |
| layer-wise lr decay [2, 9] | 0.85 (L) |
| batch size | 1024 |
| learning rate schedule | cosine decay |
| warmup epochs | 5 |
| training epochs | 50 (L) |
| augmentation | RandAug (9, 0.5) [10] |
| label smoothing [58] | 0.1 |
| erasing [74] | 0.25 |
| mixup [71] | 0.8 |
| cutmix [69] | 1.0 |
| drop path [28] | 0.1 (L) |

Table 12. ImageNet classification end-to-end fine-tuning setting.

| config | value |
| --- | --- |
| optimizer | AdamW |
| caption head lr | 1e-4 |
| other parameters lr | 1e-5 |
| weight decay | 0.01 |
| dropout | 0.1 |
| optimizer momentum | $\beta_1, \beta_2{=}0.9, 0.999$ |
| batch size | 256 |
| learning rate schedule | cosine decay |
| warmup epochs | 2 |
| training epochs | 20 |

Table 13. MSCOCO captioning end-to-end fine-tuning setting

image features from CausalFusion or DiT as input. We evaluate image features from the 14th, 21st, and 24th layers of CausalFusion and DiT, selecting the layer that achieves the highest performance. The models are fine-tuned on the Karpathy training split for 20 epochs.

## D. Additional Samples

We show more zero-shot editing results from our CausalDiffusion models in Figure 7 and 8. The editing results are achieved by first generating the original image using the initial class label, then masking a portion of the image, and regenerating it conditioned on the unmasked region and the new class label. For example, in the first example in Figure 7, an image of "volcano" is first generated. Then, the outer region of the image is masked out, and the new images are regenerated with new labels, such as "televison", "sliding door", and "car mirror".

We further show *uncurated* samples from our CausalDiffusion-XL models at 512×512 and 256×256 resolution. Figures 17 through 22 display samples under varying classifier-free guidance scales and class labels.
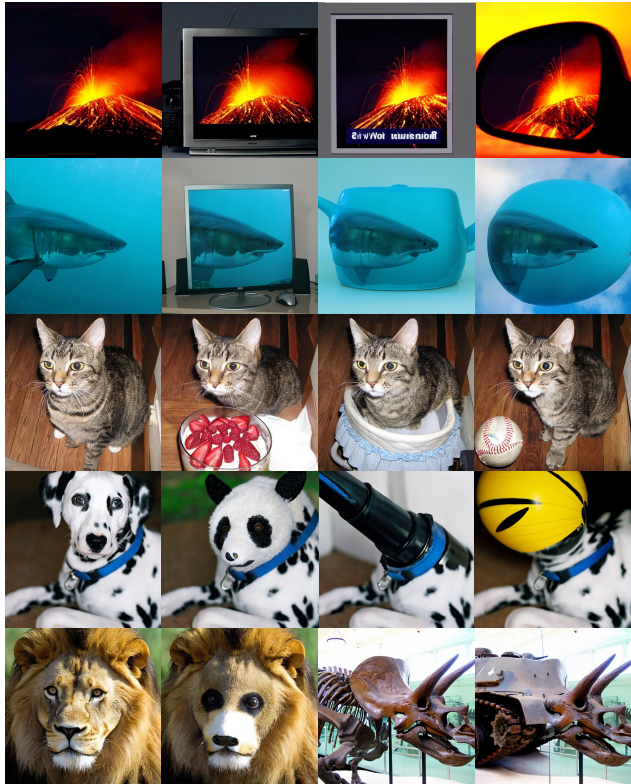


Figure 7. **Zero-shot editing samples.** CausalFusion-XL, resolution 512×512, 800 epoch, Classifier-free guidance scale = 3.0.
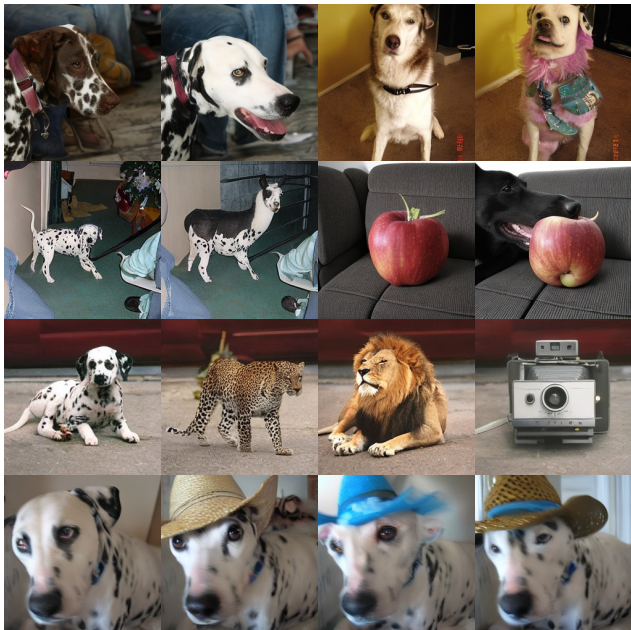


Figure 8. **Zero-shot editing samples.** CausalFusion-XL, resolution 256×256, 800 epoch, Classifier-free guidance scale = 1.5.

11

Figure 9. **Uncurated** $512 \times 512$ **CausalFusion-XL samples.**
Classifier-free guidance scale = 4.0
Class label = "otter" (360)



Figure 10. **Uncurated** $512 \times 512$ **CausalFusion-XL samples.**
Classifier-free guidance scale = 4.0
Class label = "red panda" (387)

Figure 11. **Uncurated** $512 \times 512$ **CausalFusion-XL samples.**
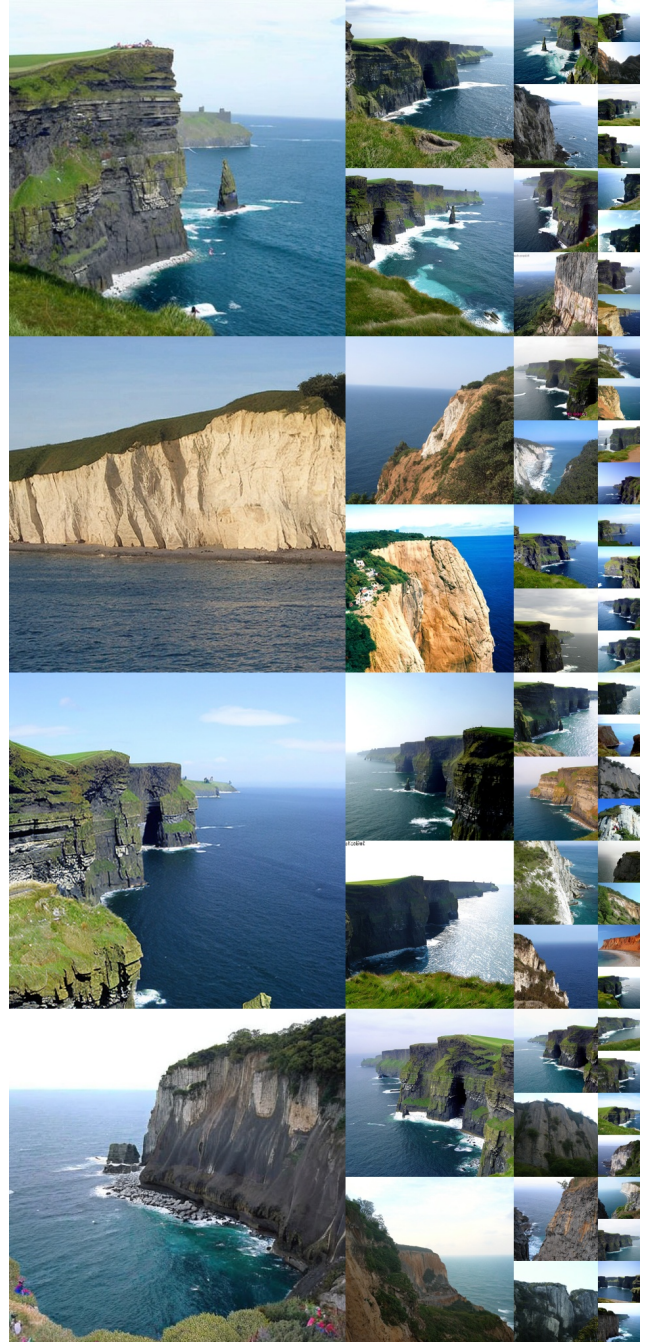Classifier-free guidance scale = 4.0
Class label = "sports car" (817)



Figure 12. **Uncurated** $512 \times 512$ **CausalFusion-XL samples.**
Classifier-free guidance scale = 4.0
Class label = "cliff" (972)

Figure 13. **Uncurated** $512 \times 512$ **CausalFusion-XL samples.**
Classifier-free guidance scale = 4.0
Class label = "arctic fox" (279)



Figure 14. **Uncurated** $512 \times 512$ **CausalFusion-XL samples.**
Classifier-free guidance scale = 4.0
Class label = "lakeshore" (975)

14

Figure 15. **Uncurated** $512 \times 512$ **CausalFusion-XL samples.**
Classifier-free guidance scale = 4.0
Class label = "sulphur-crested cockatoo" (89)



Figure 16. **Uncurated** $512 \times 512$ **CausalFusion-XL samples.**
Classifier-free guidance scale = 4.0
Class label = "geyser" (974)

Figure 17. **Uncurated** $256 \times 256$ **CausalFusion-XL samples.**
Classifier-free guidance scale = 4.0
Class label = "macaw" (88)



Figure 18. **Uncurated** $256 \times 256$ **CausalFusion-XL samples.**
Classifier-free guidance scale = 4.0
Class label = "volcano" (980)

Figure 19. **Uncurated** $256 \times 256$ **CausalFusion-XL samples.**
Classifier-free guidance scale = 2.0
Class label = "giant panda" (388)



Figure 20. **Uncurated** $256 \times 256$ **CausalFusion-XL samples.**
Classifier-free guidance scale = 2.0
Class label = "valley" (979)

Figure 21. **Uncurated** $256 \times 256$ **CausalFusion-XL samples.**
Classifier-free guidance scale = 1.5
Class label = "ice cream" (928)



Figure 22. **Uncurated** $256 \times 256$ **CausalFusion-XL samples.**
Classifier-free guidance scale = 1.5
Class label = "coral reef" (973)

# References

[1] Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth words: A vit backbone for diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22669–22679, 2023. 3, 7

[2] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. In *International Conference on Learning Representations*, 2021. 11

[3] Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. *Advances in neural information processing systems*, 28, 2015. 2

[4] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, and Eric Luhman. Video generation models as world simulators. *OpenAI Blog*, 2024. 1

[5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D 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 Ziegler, Jeffrey Wu, Clemens Winter, Chris 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. *Advances in neural information processing systems*, 33:1877–1901, 2020. 1, 2

[6] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11315–11325, 2022. 3

[7] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, José Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. In *Proceedings of the 40th International Conference on Machine Learning*, pages 4055–4075, 2023. 3

[8] Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart-\sigma: Weak-to-strong training of diffusion transformer for 4k text-to-image generation. *arXiv preprint arXiv:2403.04692*, 2024. 3

[9] Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. Electra: Pre-training text encoders as discriminators rather than generators. In *International Conference on Learning Representations*, 2020. 11

[10] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 702–703, 2020. 11

[11] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Peter Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al. Scaling vision transformers to 22 billion parameters. In *International Conference on Machine Learning*, pages 7480–7512. PMLR, 2023. 5

[12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 4

[13] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 1, 3, 4, 7, 10

[14] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, and Jie Tang. Cogview: Mastering text-to-image generation via transformers. *Advances in neural information processing systems*, 34:19822–19835, 2021. 3

[15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 10

[16] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv:2407.21783*, 2024. 1, 2, 8

[17] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024. 1

[18] Zhengcong Fei, Mingyuan Fan, Changqian Yu, Debang Li, and Junshi Huang. Scaling diffusion transformers to 16 billion parameters. *arXiv preprint arXiv:2407.11633*, 2024. 7

[19] 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. In *European Conference on Computer Vision*, pages 89–106. Springer, 2022. 3

[20] Peng Gao, Le Zhuo, Ziyi Lin, Chris Liu, Junsong Chen, Ruoyi Du, Enze Xie, Xu Luo, Longtian Qiu, Yuhang Zhang, et al. Lumina-t2x: Transforming text into any modality, resolution, and duration via flow-based large diffusion transformers. *arXiv preprint arXiv:2405.05945*, 2024. 7

[21] Jiatao Gu, Yuyang Wang, Yizhe Zhang, Qihang Zhang, Dinghuai Zhang, Navdeep Jaitly, Josh Susskind, and Shuangfei Zhai. Dart: Denoising autoregressive transformer for scalable text-to-image generation. *arXiv preprint arXiv:2410.08159*, 2024. 1, 3, 7

[22] Tiankai Hang, Shuyang Gu, Chen Li, Jianmin Bao, Dong Chen, Han Hu, Xin Geng, and Baining Guo. Efficient diffusion training via min-snr weighting strategy. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7441–7451, 2023. 2, 7

[23] Shaozhe Hao, Xuantong Liu, Xianbiao Qi, Shihao Zhao, Bojia Zi, Rong Xiao, Kai Han, and Kwan-Yee K Wong. Bigr: Harnessing binary latent codes for image generation and

improved visual representation capabilities. *arXiv preprint arXiv:2410.14672*, 2024. 3, 7

[24] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022. 4, 8, 10

[25] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 7, 8

[26] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1, 2, 3, 10

[27] Emiel Hoogeboom, Jonathan Heek, and Tim Salimans. simple diffusion: End-to-end diffusion for high resolution images. In *International Conference on Machine Learning*, pages 13213–13232. PMLR, 2023. 7

[28] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 646–661. Springer, 2016. 11

[29] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Advances in neural information processing systems*, 35:26565–26577, 2022. 1

[30] Diederik Kingma and Ruiqi Gao. Understanding diffusion objectives as the elbo with simple data augmentation. *Advances in Neural Information Processing Systems*, 36, 2024. 7

[31] Dan Kondratyuk, Lijun Yu, Xiuye Gu, Jose Lezama, Jonathan Huang, Grant Schindler, Rachel Hornung, Vighnesh Birodkar, Jimmy Yan, Ming-Chang Chiu, et al. Videopoet: A large language model for zero-shot video generation. In *Forty-first International Conference on Machine Learning*, 2024. 3

[32] Tuomas Kynkäänniemi, Miika Aittala, Tero Karras, Samuli Laine, Timo Aila, and Jaakko Lehtinen. Applying guidance in a limited interval improves sample and distribution quality in diffusion models. *arXiv preprint arXiv:2404.07724*, 2024. 7

[33] Black Forest Labs. Flux, 2024. 3

[34] Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image generation without vector quantization. *arXiv preprint arXiv:2406.11838*, 2024. 3, 7, 8, 9, 10

[35] Yanghao Li, Haoqi Fan, Ronghang Hu, Christoph Feichtenhofer, and Kaiming He. Scaling language-image pre-training via masking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23390–23400, 2023. 4, 10

[36] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 8

[37] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2023. 1

[38] Qihao Liu, Zhanpeng Zeng, Ju He, Qihang Yu, Xiaohui Shen, and Liang-Chieh Chen. Alleviating distortion in image generation via multi-resolution diffusion models. *arXiv preprint arXiv:2406.09416*, 2024. 7

[39] I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 10

[40] Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha Kembhavi. Unified-io 2: Scaling autoregressive multimodal models with vision language audio and action. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26439–26455, 2024. 3

[41] Zhuoyan Luo, Fengyuan Shi, Yixiao Ge, Yujiu Yang, Limin Wang, and Ying Shan. Open-magvit2: An open-source project toward democratizing auto-regressive visual generation. *arXiv preprint arXiv:2409.04410*, 2024. 7

[42] Nanye Ma, Mark Goldstein, Michael S Albergo, Nicholas M Boffi, Eric Vanden-Eijnden, and Saining Xie. Sit: Exploring flow and diffusion-based generative models with scalable interpolant transformers. *arXiv preprint arXiv:2401.08740*, 2024. 7

[43] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 3

[44] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023. 1, 2, 3, 4, 5, 6, 7, 8, 9

[45] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. In *The Twelfth International Conference on Learning Representations*, 2024. 3

[46] Alec Radford and Karthik Narasimhan. Improving language understanding by generative pre-training. *OpenAI blog*, 2018. 1, 2, 8

[47] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 2019. 1, 2

[48] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pages 8821–8831. Pmlr, 2021. 3

[49] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv:2204.06125*, 2022. 3

[50] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1, 3, 4, 7

[51] 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. 7

[52] 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. 3

[53] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019. 3

[54] 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

[55] StabilityAI. https://huggingface.co/stabilityai/sd-vae-ft-ema, 2024. 10

[56] Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan. Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint arXiv:2406.06525*, 2024. 7

[57] Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiying Yu, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative multimodal models are in-context learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14398–14409, 2024. 3

[58] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016. 11

[59] Jiayan Teng, Wendi Zheng, Ming Ding, Wenyi Hong, Jianqiao Wangni, Zhuoyi Yang, and Jie Tang. Relay diffusion: Unifying diffusion process across resolutions for image synthesis. *arXiv preprint arXiv:2309.03350*, 2023. 7

[60] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive modeling: Scalable image generation via next-scale prediction. *arXiv preprint arXiv:2404.02905*, 2024. 7

[61] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv:2302.13971*, 2023. 1, 2

[62] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv:2307.09288*, 2023. 1, 2

[63] Michael Tschannen, Cian Eastwood, and Fabian Mentzer. Givt: Generative infinite-vocabulary transformers. In *European Conference on Computer Vision*, pages 292–309. Springer, 2025. 7

[64] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv:2409.12191*, 2024. 8, 10

[65] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018. 6

[66] ZiDong Wang, Zeyu Lu, Di Huang, Cai Zhou, Wanli Ouyang, et al. Fitv2: Scalable and improved flexible vision transformer for diffusion model. *arXiv preprint arXiv:2410.13925*, 2024. 7

[67] Chen Wei, Haoqi Fan, Saining Xie, Chao-Yuan Wu, Alan Yuille, and Christoph Feichtenhofer. Masked feature prediction for self-supervised visual pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14668–14678, 2022. 4, 8

[68] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *Transactions on Machine Learning Research*, 2022. 3

[69] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6023–6032, 2019. 11

[70] Jun Zhan, Junqi Dai, Jiasheng Ye, Yunhua Zhou, Dong Zhang, Zhigeng Liu, Xin Zhang, Ruibin Yuan, Ge Zhang, Linyang Li, et al. Anygpt: Unified multimodal llm with discrete sequence modeling. *arXiv:2402.12226*, 2024. 3

[71] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018. 11

[72] Chuyang Zhao, Yuxing Song, Wenhao Wang, Haocheng Feng, Errui Ding, Yifan Sun, Xinyan Xiao, and Jingdong Wang. Monoformer: One transformer for both diffusion and autoregression. *arXiv:2409.16280*, 2024. 1, 7

[73] Hongkai Zheng, Weili Nie, Arash Vahdat, and Anima Anandkumar. Fast training of diffusion models with masked transformers. *arXiv preprint arXiv:2306.09305*, 2023. 7

[74] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *Proceedings of the AAAI conference on artificial intelligence*, pages 13001–13008, 2020. 11

[75] Chunting Zhou, Lili Yu, Arun Babu, Kushal Tirumala, Michihiro Yasunaga, Leonid Shamis, Jacob Kahn, Xuezhe Ma, Luke Zettlemoyer, and Omer Levy. Transfusion: Predict the next token and diffuse images with one multi-modal model. *arXiv:2408.11039*, 2024. 1, 3, 8