# AUTOREGRESSIVE VIDEO GENERATION WITHOUT VECTOR QUANTIZATION

**Haoge Deng**[1,4*], **Ting Pan**[2,4*], **Haiwen Diao**[3,4*], **Zhengxiong Luo**[4*], **Yufeng Cui**[4], **Huchuan Lu**[3], **Shiguang Shan**[2], **Yonggang Qi**[1], **Xinlong Wang**[4†]
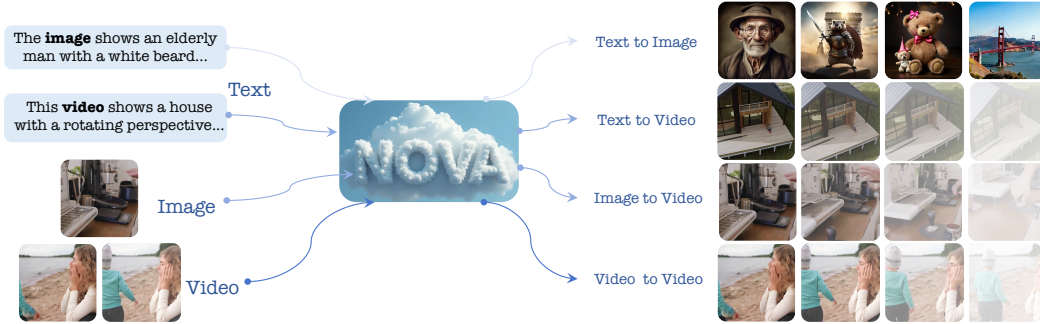
BUPT[1], ICT-CAS[2], DLUT[3], BAAI[4]

Figure 1: **NOVA** is a non-quantized autoregressive model for efficient and flexible visual generation.

## ABSTRACT

This paper presents a novel approach that enables autoregressive video generation with high efficiency. We propose to reformulate the video generation problem as a non-quantized autoregressive modeling of temporal *frame-by-frame* prediction and spatial *set-by-set* prediction. Unlike raster-scan prediction in prior autoregressive models or joint distribution modeling of fixed-length tokens in diffusion models, our approach maintains the causal property of GPT-style models for flexible in-context capabilities, while leveraging bidirectional modeling within individual frames for efficiency. With the proposed approach, we train a novel video autoregressive model without vector quantization, termed **NOVA**. Our results demonstrate that **NOVA** surpasses prior autoregressive video models in data efficiency, inference speed, visual fidelity, and video fluency, even with a much smaller model capacity, *i.e.*, 0.6B parameters. **NOVA** also outperforms state-of-the-art image diffusion models in text-to-image generation tasks, with a significantly lower training cost. Additionally, **NOVA** generalizes well across extended video durations and enables diverse zero-shot applications in one unified model. Code and models are publicly available at https://github.com/baaivision/NOVA.

## 1 INTRODUCTION

Autoregressive large language models (LLMs) (Brown et al. (2020); Touvron et al. (2023)) have become a foundational architecture in natural language processing (NLP), exhibiting emerging capabilities in in-context learning and long-context reasoning. In autoregressive (AR) vision generation domain, prior approaches (Ramesh et al. (2021); Ding et al. (2021); Yu et al. (2022); Yan et al. (2021); Villegas et al. (2022); Kondratyuk et al. (2023); Wang et al. (2024)) typically transform images or video clips into a discrete-valued token space using vector quantization (Van Den Oord et al. (2017); Esser et al. (2021)), which are then flattened into sequences for token-by-token prediction. However, it is challenging for vector-quantized tokenizers to achieve high fidelity and high compression simultaneously. More tokens are required for high quality. Thus, the cost increases substantially with higher image resolutions or longer video sequences.

---

*Equal Contribution,†Corresponding Author: *wangxinlong@baai.ac.cn*

In contrast, video diffusion models (Brooks et al. (2024); Kuaishou (2024); Blattmann et al. (2023)) learn with highly compressed video sequences in a compact continuous latent space. However, most of them only learn the joint distribution of fixed-length frames, lacking the flexibility to generate videos with varied lengths. More importantly, they do not possess the in-context abilities of autoregressive models, *i.e.*, solving diverse tasks in context with a unified model such as GPT for language.

In this work, we present **NOVA**, which addresses the issues above and enables autoregressive video generation with high efficiency. We propose to reformulate the video generation problem as a non-quantized autoregressive modeling of temporal *frame-by-frame* prediction and spatial *set-by-set* prediction. **NOVA** is inspired by Emu3 (Wang et al. (2024)) for autoregressive video and multimodal generation, and MAR (Li et al. (2024c)) for non-quantized autoregressive image generation, which utilizes non-quantized vectors as visual tokens and performs set-by-set autoregressive prediction. While both are non-quantized autoregressive approaches, it is non-trivial at all from MAR to **NOVA**: **1) NOVA** solves the challenges including efficiency, scalability, and mask schedule when learning more complex text-to-image generation instead of class-to-image generation. **2) NOVA** first predicts temporal frames sequentially and then predicts spatial sets within each frame. **NOVA** is the first to enable a non-quantized autoregressive model for video generation.

Specifically, **NOVA** predicts each frame in a casual order temporally, and predicts each token set in a random order spatially. In this way, text-to-video generation can be regarded as a fundamental task that implicitly and comprehensively encompasses various generation tasks (See Figure 1), including text-to-image, image-to-video, text&image-to-video, *etc*. With non-quantized tokenizers and a flexible autoregressive framework, **NOVA** simultaneously takes advantage of **1**) high-fidelity and compact visual compression for low cost in training and inference, and **2**) in-context abilities for integrating multiple visual generation tasks in a unified model.

For text-to-video generation, **NOVA** surpasses autoregressive counterparts in data efficiency, inference speed, and video fluency, while matching the performance of diffusion models of similar scale, *e.g.*, achieving a VBench (Huang et al. (2024)) score of 80.1 with a processing speed of 2.75 FPS[1], trained in only 342 GPU days on A100-40G. For text-to-image generation, NOVA achieves a GenEval (Ghosh et al. (2024)) score of 0.75, surpassing previous diffusion models with notably lower training cost, *e.g.*, only 127 GPU days for training this state-of-the-art 0.6B model. Additionally, **NOVA** also demonstrates strong zero-shot generalization across various contexts. We believe that **NOVA** paves the way for next-generation video generation, offering possibilities for real-time and infinite video generation, beyond Sora-like video diffusion models.

## 2 RELATED WORKS

### 2.1 DIFFUSION MODELS FOR VISUAL GENERATION

Diffusion models (Ho et al. (2020); Song et al. (2020)) have made significant advances in visual generation, including text-to-image tasks (Esser et al. (2024a); Betker et al. (2023a); Baldridge et al. (2024)) and text-to-video tasks (Brooks et al. (2024); Lin et al. (2024); Blattmann et al. (2023)). Image diffusion models typically model the joint distribution of fixed-length tokens in pixel (Ho et al. (2020); Nichol et al. (2021); Hoogeboom et al. (2023)) or latent space (Rombach et al. (2022a); Esser et al. (2024a); Betker et al. (2023a); Chen et al. (2023)). Besides, video diffusion models further introduce temporal layers to capture relationships between a fixed number of video frames. After training, additional tasks and modalities are added by incorporating extra inference tricks (Meng et al. (2021)), structure moderation (Blattmann et al. (2023); Esser et al. (2023); Liew et al. (2023)), and adapter layers (Zhang et al. (2023b); Guo et al. (2023)). Although these strategies can be composable, they stand in contrast to the autoregressive approaches (Kondratyuk et al. (2023); Hong et al. (2022); Radford (2018); Touvron et al. (2023)), which trains a single model end-to-end for multi-task learning, offering notable context scalability and zero-shot generalizability across diverse application scenarios, especially in extending video generation duration.

---

[1]The 2.75 FPS is measured on a single NVIDIA A100-40G GPU using a batch size of 24.

## 2.2 Autoregressive Models for Visual Generation

**Raster-scan Autoregressive Models** are typically implemented on the discrete-valued RGB pixels (Kalchbrenner et al. (2017); Reed et al. (2017)) or latent space (Esser et al. (2021); Van Den Oord et al. (2017)), analogous to their language counterparts (Radford et al. (2019); Anil et al. (2023)). Recent studies involve scalable autoregressive transformers to generate token sequences in the raster-scan order for image generation (Ramesh et al. (2021); Ding et al. (2021; 2022); Yu et al. (2022); Sun et al. (2024b)), and video generation (Yan et al. (2021); Kondratyuk et al. (2023); Nash et al. (2022)). Specifically, VAR (Tian et al. (2024)) introduces next-scale prediction to progressively process the token-by-token sequence across multiple resolutions, leading to improved image quality.

**Masked Autoregressive Models** further develop a masked generative models (Chang et al. (2022)) to introduce a generalized autoregressive concept. They introduce a bidirectional transformer and predict randomly masked tokens by attending to unmasked conditions. This makes up for the suboptimal modeling and inefficient inference of sequentially line-by-line strategy, which inspires a series of subsequent works in text-to-image (Chang et al. (2023)) and text-to-video generation (Hong et al. (2022); Yu et al. (2023); Villegas et al. (2022)). Particularly, MAR (Li et al. (2024c)) decouples discrete tokenizers from autoregressive models and utilizes a diffusion procedure for per-token probability distributions. It is fully validated in the class-to-image field, holding great potential in the text-to-image domain. However, its application to text-to-video generation intuitively requires a masked autoregressive process across entire video frames, challenging multi-context learning and training efficiency. In contrast, our NOVA model breaks down video generation into frame-by-frame temporal predictions combined with spatial set-by-set predictions. This allows each frame to act as a meta causal unit, enabling extended video duration and zero-shot generalizability across various contexts. Besides, the subsequent spatial set-of-tokens prediction unlocks the power of bidirectional modeling patterns, enhancing inference efficiency while preserving visual quality and fidelity.

## 3 Methodology

We first review two categories of autoregressive video generation in Sec. 3.1. In Sec. 3.2-3.4, we introduce framework pipeline and implementation details of our NOVA, illustrated in Figure 2.

## 3.1 Rethinking autoregressive models for video generation

As mentioned above, we regard text-to-video generation and autoregressive (AR) model as the basic task and means, respectively. We briefly retrospect related technical background. There exist two types of AR video generation approaches: **(1) Token-by-token generation via raster-scan order.** These studies perform causal per-token prediction within video frame sequence (Kondratyuk et al. (2023)), and decode vision tokens sequentially following the raster scan ordering (Wang et al. (2024)), which is defined as follows:

Table 1: Symbology Settings.

| | |
|---|---|
| $N, n$ | The number of all video tokens. |
| $F, f$ | The number of all video frames. |
| $K, k$ | The number of sets in an image. |

$$p\left(C, x_1, ..., x_N\right) = \prod_{n}^{N} p\left(x_n \mid C, x_1, ..., x_{n-1}\right), \tag{1}$$

where $C$ indicates various condition contexts, *e.g.*, label, text, image, and *etc*. Note that $x_n$ denotes $n$-th token of $N$ video raster-scale tokens. In contrast, **(2) Masked set-by-set generation in a random order** treats all tokens within each video frame equally, using a bidirectional transformer decoder for per-set prediction (Yu et al. (2023)). However, this generalized autoregressive (AR) model is trained using synchronous modeling on large, fixed-length video frames, which can lead to poor scalability in context and issues with coherence over longer video durations. Hence, NOVA proposes a novel solution by decoupling per-set generation within a single video frame from the per-frame prediction across the entire video sequence. This allows NOVA to better handle both temporal causality and spatial relationships, providing a more flexible and scalable AR framework.
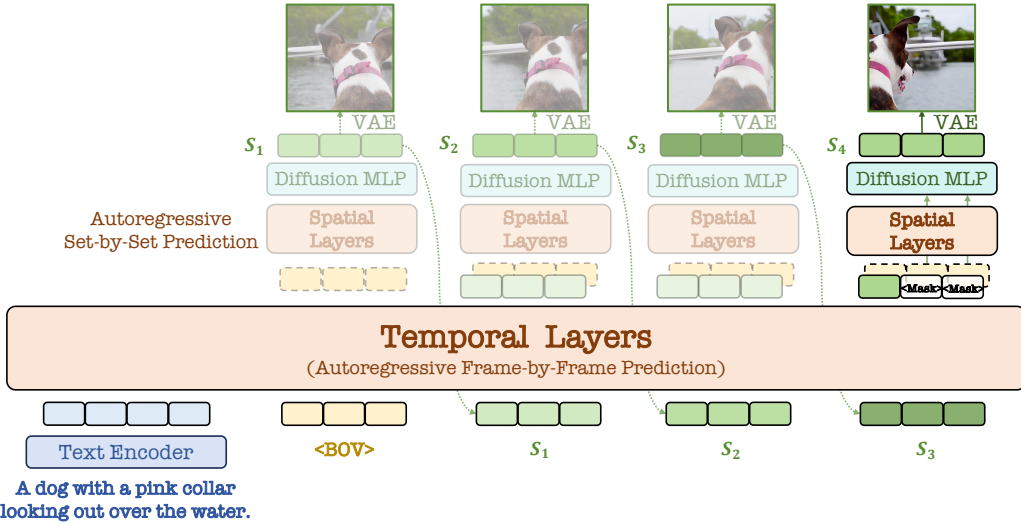
Figure 2: **NOVA framework and the inference process.** With text inputs, NOVA performs autoregressive generation via temporal frame-by-frame prediction and spatial set-by-set prediction. Finally, we implement diffusion denoising in a continuous-values space.

## 3.2 TEMPORAL AUTOREGRESSIVE MODELING VIA FRAME-BY-FRAME PREDICTION

Inspired by (Zhuo et al. (2024)), we use a pre-trained language model (Javaheripi et al. (2023)) to encode text prompts to features. To better control video dynamics, we use OpenCV (cv2) (Bradski (2000)) to compute the optical flow of sampled video frames. The average flow magnitude is used as a motion score and integrated with the prompt. Besides, we employ open-source 3D variational autoencoder (VAE) (Lin et al. (2024)) with a temporal stride of 4 and a spatial stride of 8 to encode the video frames to the latent space. We add an additional learnable patch embedding layer with a spatial stride of 4 to align channels of latent video to the subsequent transformer. Notably, next-token prediction in early AR models seems counter-intuitive for undirected visual patches within a single image and suffers from high latency during inference. In contrast, video frames can naturally be viewed as a causal sequence, with each frame acting as a meta unit for AR generation. Therefore, we implement block-wise causal masking attention depicted in Figure 3(a), ensuring that each frame can only attend to the text prompts, video flow, and its preceding frames, while allowing all current frame tokens to be visible to each other as follows:

$$p\left(P, m, B, S_1, ..., S_F\right) = \prod_f^F p\left(S_f \mid P, m, B, S_1, ..., S_{f-1}\right), \quad (2)$$

where $P, m$ indicate text prompts and video flow respectively. Here, $S_f$ denotes the overall tokens of $f$-th video frame, and $B$ represent learnable begin-of-video (BOV) embeddings for predicting the initial video frame, the number of which corresponds to the patch number of one single frame. Note that we add 1-D and 2-D sine-cosine embeddings (Vaswani et al. (2017)) with video frame features to indicate time and position information respectively, which are convenient for temporal and spatial extrapolation. From equation 2, we can reformulate text-to-image and image-to-video generation as $p\left(S_1 \mid P, m, B\right)$ and $p\left(S_f \mid \varnothing, m, B, S_1, ..., S_{f-1}\right)$. This generalized causal process can synchronously model the condition contexts for each video frame, greatly enhancing training efficiency, and allowing the kv-cache technology for fast decoding procedure during inference.

## 3.3 SPATIAL AUTOREGRESSIVE MODELING VIA SET-BY-SET PREDICTION

Inspired by (Chang et al. (2022); Li et al. (2024c)), we define each token set with multiple tokens from random directions as a meta causal token in Figure 3(b), facilitating a generalized AR process with efficient parallel decoding. Notably, we tried to utilize the temporal layers' outputs targeting one frame as indicator features to assist the spatial layers, gradually decoding all randomly masked token sets within the corresponding image. However, this approach resulted in image structure collapse and
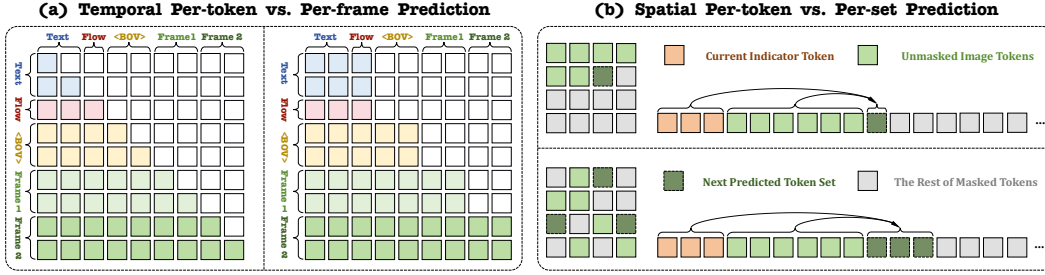
4

Figure 3: Overview of our block-wise temporal and spatial generalized autoregressive attention. Different from per-token generation, NOVA regressively predicts each frame in a casual order across the temporal scale, and predicts each token set in a random order across the spatial scale.

inconsistent video fluency over the increasing number of frames. We hypothesize that this occurs because the indicator features from adjacent frames are similar, making it difficult to accurately learn continuous and imperceptible motion changes without explicit modeling. Besides, the indicator features derived from the ground-truth contextual frame during training contribute to weak robustness and stability of spatial AR layers against cumulative inference errors.

To address this issue, we introduce a Scaling and Shift Layer that reformulates cross-frame motion changes by learning relative distribution variations within a unified space, rather than directly modeling the unreferenced distribution of the current frame. Notably, we select the BOV-attended output of the temporal layers as the anchor feature set, as it serves as the initial feature set with significantly less noise accumulation than subsequent frame feature sets. Specifically, we first translate the features from current frame set into dimension-wise variance and mean parameters $\gamma$ and $\beta$ via multi-layer perception (MLP). After that, we affine the normalized features from the anchor set into indicator features $S_f'$ via channel-wise scale and shift operation. Specially, we explicitly set $\gamma = 1$ and $\beta = 0$ for the first frame. With unmasked token features, we predict randomly masked visual tokens in a set-by-set order through a bidirectional paradigm, which can be formulated as follows:

$$p\left(S_f', S_{(f,1)}, ..., S_{(f,K)}\right) = \prod_k^K p\left(S_{(f,k)} \mid S_f', S_{(f,1)}, ..., S_{(f,k-1)}\right), \tag{3}$$

where $S_f'$ denotes the indicator features for generating $f$-th video frame, and $S_{(f,k)}$ denotes $k$-th token set of $f$-th video frame. We add 2-D sine-cosine embeddings with masked and unmasked tokens to indicate their relative position. This generalized spatial AR prediction leverages powerful bidirectional patterns within single-image tokens and achieves efficient inference with parallel masked decoding. *Notably, we incorporate post-norm layers before the residual connections in both temporal and spatial AR layers.* Our empirical findings show that this design effectively addresses architectural and optimization challenges that previously hindered stable training in generalized video generation.

### 3.4 DIFFUSION PROCEDURE DENOISING FOR PER-TOKEN PREDICTION

During training, we import *diffusion loss* (Li et al. (2024c)) to estimate per-token probability in a continuous-valued space. For example, we define one ground-truth token as $x_n$ and corresponding NOVA's output as $z_n$. The loss function can be formulated as a denoising criterion:

$$\mathcal{L}(x_n \mid z_n) = \mathbb{E}_{\varepsilon,t}\left[\left\|\epsilon - \epsilon_\theta\left(x_n^t \mid t, z_n\right)\right\|^2\right]. \tag{4}$$

Here $\epsilon$ is a Gaussian vector sampled from $\mathcal{N}(\mathbf{0}, \mathbf{I})$, and noisy data $x_n^t = \sqrt{\bar{\alpha}_t}x_n + \sqrt{1 - \bar{\alpha}_t}\epsilon$, where $\bar{\alpha}_t$ is a noise schedule (Nichol & Dhariwal (2021)) indexed by a time step $t$. The noise estimator $\epsilon_\theta$ is multiple MLP blocks parameterized by $\theta$. The notation $\epsilon_\theta(x_n^t \mid t, z_n)$ means that this network takes $x_n^t$ as the input, and is conditional on both $t$ and $z_n$. We follow (Li et al. (2024c)) to sample $t$ by 4 times during training for each image.

During inference, we sample $x_n^T$ from a random Gaussian noise $\mathcal{N}(\mathbf{0}, \mathbf{I})$ and denoise it step-by-step by sequentially sampling $x_n^T$ to $x_n^0$ via $x_n^{t-1} = \frac{1}{\sqrt{\alpha_t}}\left(x_n^t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon_\theta\left(x_n^t|t,z_n\right)\right) + \sigma_t\epsilon$, where $\sigma_t$ is the noise level at time step $t$, and $\epsilon$ is sampled from the Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

# 4 EXPERIMENT

## 4.1 EXPERIMENT SETUP

**Datasets.** We involve several diverse, curated, and high-quality datasets to facilitate the training of our NOVA. For text-to-image training, we initially curate 16M image-text pairs sourced from DataComp (Gadre et al. (2024)), COYO (Byeon et al. (2022)), Unsplash (UnsplashTeam (2020)), and JourneyDB (Sun et al. (2024a)). To explore the scaling properties of NOVA, we expanded the dataset to approximately 600M image-text pairs by selecting more images that have a minimum aesthetic score of 5.0 from LAION (Schuhmann et al. (2022)), DataComp and COYO. For text-to-video training, we select 19M video-text pairs on a subset (Lin et al. (2024)) of Panda-70M (Chen et al. (2024c)) and internal video-text pairs. We further collect 1M of high-resolution video-text pairs from Pexels (PexelsTeam (2014)) to fine-tune our final video generation model. Following (Diao et al. (2024)), we train a caption engine based on Emu2-17B (Sun et al. (2023)) model to create high-quality descriptions for our image and video datasets. The maximum text length is set to 256.

**Architectures.** We mostly follow (Li et al. (2024c)) to build NOVA's spatial AR layer and denoising MLP block, including a layer sequence of LayerNorm (Lei Ba et al. (2016)), AdaLN (Huang & Belongie (2017a)), linear layer, SiLU activation (Elfwing et al. (2018)), and another linear layer. We configure the temporal encoder, spatial encoder, and decoder with 16 layers each, using a dimension of 768 (0.3B), 1024 (0.6B) or 1536 (1.4B). The denoising MLP consists of 3 blocks with a dimension of 1280. The spatial layers adopt the encoder-decoder architecture of MAR (Li et al. (2024c)), similar to MAE (He et al. (2022)). Specifically, the encoder processes the visible patches for reconstruction. The decoder further processes visible and masked patches for generation. To capture the image latent features, we employ a pre-trained and frozen VAE from (Lin et al. (2024)), which achieves $4\times$ compression in the temporal dimension and $8 \times 8$ compression in the spatial dimension. We adopt the masking and diffusion schedulers from (Li et al. (2024c); Nichol & Dhariwal (2021)), using a masking ratio between 0.7 and 1.0 during training, and progressively reducing it from 1.0 to 0 following a cosine schedule (Chang et al. (2023)) during inference. In line with common practice (Ho et al. (2020)), we train with a 1000-step noise schedule but default to 100 steps for inference.

**Training details.** NOVA is trained with sixteen A100 (40G) nodes. We utilize the AdamW optimizer (Loshchilov et al. (2017)) ($\beta_1 = 0.9, \beta_2 = 0.95$) with a weight decay of 0.02 and a base learning rate of 1e-4 in all experiments. The peak learning rate is adjusted for different batch sizes during training using the scaling rule (Goyal (2017)) : lr = base_lr $\times$ batchsize$/256$. We train text-to-image models from scratch and then load these weights to train text-to-video models.

**Evaluation.** We use T2I-CompBench (Huang et al. (2023)), GenEval (Ghosh et al. (2024)) and DPG-Bench (Hu et al. (2024)) to assess the alignment between the generated images and text condition. We generate image samples for each of the original or rewritten (Wang et al. (2024)) text prompts. Each image sample has a resolution of 512×512 or 1024×1024. We use VBench (Huang et al. (2024)) to evaluate the capacity of text-to-video generation across 16 dimensions. For a given text prompt, we randomly generate 5 samples, each with a video size of 33×768×480. We employ classifier-free guidance (Ho & Salimans (2022)) with a value of 7.0 along with 128 autoregressive steps to enhance the quality of the generated images and videos in all evaluation experiments.

## 4.2 MAIN RESULTS

**NOVA outperforms existing text-to-image models with superior performance and efficiency.** In Table 2, we compare NOVA with several recent text-to-image models, including PixArt-$\alpha$ (Chen et al. (2023)), SD v1/v2 (Rombach et al. (2022b)), SDXL (Podell et al. (2023)), DALL-E2 (Ramesh et al. (2022)), DALL-E3 (Betker et al. (2023b)), SD3 (Esser et al. (2024b)), LlamaGen (Sun et al. (2024b) and Emu3 (Wang et al. (2024)). After text-to-image training, NOVA achieves state-of-the-art performance on the GenEval benchmark, especially in generating a specified number of targets. Notably, NOVA also achieves leading results on T2I-CompBench and DPG-Bench, excelling at both the small model scale and data scale (requiring only 16% training overhead of the best competitor PixArt-$\alpha$). *Last but not least, our text-to-video model outperforms most specialized text-to-image models, e.g., SD v1/v2, SDXL and DALL-E2.* This underscores the robustness and versatility of our model in multi-context scenarios, with text-to-video generation as the fundamental training task.

Table 2: **Text-to-image evaluation on various benchmarks.** The best and second-best results are in blue and green . The data is from Huang et al. (2023),Wang et al. (2024) and Esser et al. (2024b).

| Model | ModelSpec | | T2I-CompBench | | | GenEval | | | | | | | DPG-Bench | A100 days |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | #params | #images | Color | Shape | Texture | Overall | Single | Two | Counting | Colors | Position | ColorAttr | Overall | |
| *Diffusion models* | | | | | | | | | | | | | | |
| PixArt-$\alpha$ | 0.6B | 25M | 68.86 | 55.82 | 70.44 | 0.48 | 0.98 | 0.50 | 0.44 | 0.80 | 0.08 | 0.07 | 71.11 | 753 |
| SD v1.5 | 1B | 2B | 37.50 | 37.24 | 42.19 | 0.43 | 0.97 | 0.38 | 0.35 | 0.76 | 0.04 | 0.06 | 63.18 | - |
| SD v2.1 | 1B | 2B | 56.94 | 44.95 | 49.82 | 0.50 | 0.98 | 0.37 | 0.44 | 0.85 | 0.07 | 0.17 | - | - |
| SDXL | 2.6B | - | 63.69 | 54.08 | 56.37 | 0.55 | 0.98 | 0.44 | 0.39 | 0.85 | 0.15 | 0.23 | 74.65 | - |
| DALL-E2 | 6.5B | 650M | 57.50 | 54.64 | 63.74 | 0.52 | 0.94 | 0.66 | 0.49 | 0.77 | 0.10 | 0.19 | - | - |
| DALL-E3 | - | - | 81.10 | 67.50 | 80.70 | 0.67 | 0.96 | 0.87 | 0.47 | 0.83 | 0.43 | 0.45 | 83.50 | - |
| SD3 | 2B | - | - | - | - | 0.62 | 0.98 | 0.74 | 0.63 | 0.67 | 0.34 | 0.36 | 84.10 | - |
| *Autoregressive models* | | | | | | | | | | | | | | |
| LlamaGen | 0.8B | 60M | - | - | - | 0.32 | 0.71 | 0.34 | 0.21 | 0.58 | 0.07 | 0.04 | - | - |
| Emu3 (+ Rewriter) | 8B | - | 79.13 | 58.46 | 74.22 | 0.66 | 0.99 | 0.81 | 0.42 | 0.80 | 0.49 | 0.45 | 81.60 | - |
| NOVA (512×512) | 0.6B | 16M | 70.75 | 55.98 | 69.79 | 0.66 | 0.98 | 0.85 | 0.58 | 0.83 | 0.20 | 0.48 | 81.76 | 127 |
| + Rewriter | 0.6B | 16M | 83.02 | 61.47 | 75.80 | 0.75 | 0.98 | 0.88 | 0.62 | 0.82 | 0.62 | 0.58 | - | 127 |
| + Videos | 0.6B | 36M | 71.80 | 47.86 | 65.31 | 0.55 | 0.98 | 0.56 | 0.48 | 0.75 | 0.15 | 0.41 | 81.77 | 342 |
| + Videos & Rewriter | 0.6B | 36M | 81.36 | 59.16 | 72.45 | 0.71 | 0.98 | 0.83 | 0.52 | 0.81 | 0.58 | 0.51 | - | 342 |
| NOVA (1024×1024) | 0.3B | 600M | 73.35 | 57.28 | 70.09 | 0.67 | 0.98 | 0.86 | 0.53 | 0.84 | 0.32 | 0.52 | 80.60 | 267 |
| NOVA (1024×1024) | 0.6B | 600M | 74.72 | 56.99 | 69.50 | 0.69 | 0.98 | 0.89 | 0.56 | 0.84 | 0.32 | 0.56 | 82.25 | 320 |
| NOVA (1024×1024) | 1.4B | 600M | 74.30 | 57.14 | 70.00 | 0.71 | 0.99 | 0.91 | 0.62 | 0.85 | 0.33 | 0.56 | 83.01 | 608 |

Table 3: **Text-to-video evaluation on VBench.** We have classified existing video generation methods into different categories for better clarity. The baseline data is sourced from Huang et al. (2024).

| Model | #params | #videos | latency | Total Score | Quality Score | Semantic Score | Aesthetic Quality | Object Class | Multiple Objects | Human Action | Spatial Relationship | Scene |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Closed-source models* | | | | | | | | | | | | |
| Gen-2 | - | - | - | 80.58 | 82.47 | 73.03 | 66.96 | 90.92 | 55.47 | 89.20 | 66.91 | 48.91 |
| Kling (2024-07) | - | - | - | 81.85 | 83.39 | 75.68 | 61.21 | 87.24 | 68.05 | 93.40 | 73.03 | 50.86 |
| Gen-3 | - | - | - | 82.32 | 84.11 | 75.17 | 63.34 | 87.81 | 53.64 | 96.4 | 65.09 | 54.57 |
| *Diffusion models (w/ SD init)* | | | | | | | | | | | | |
| LaVie | 3B | 25M | - | 77.08 | 78.78 | 70.31 | 54.94 | 91.82 | 33.32 | 96.8 | 34.09 | 52.69 |
| Show-1 | 4B | 10M | - | 78.93 | 80.42 | 72.98 | 57.35 | 93.07 | 45.47 | 95.60 | 53.50 | 47.03 |
| AnimateDiff-v2 | 1B | 10M | - | 80.27 | 82.90 | 69.75 | 67.16 | 90.90 | 36.88 | 92.60 | 34.60 | 50.19 |
| VideoCrafter-v2.0 | 2B | 10M | - | 80.44 | 82.20 | 73.42 | 63.13 | 92.55 | 40.66 | 95.00 | 35.86 | 55.29 |
| T2V-Turbo (VC2) | 2B | 10M | - | 81.01 | 82.57 | 74.76 | 63.04 | 93.96 | 54.65 | 95.20 | 38.67 | 55.58 |
| *Diffusion models* | | | | | | | | | | | | |
| OpenSora-v1.1 | 1B | 10M | 48s | 75.66 | 77.74 | 67.36 | 50.12 | 86.76 | 40.97 | 84.20 | 52.47 | 38.63 |
| OpenSoraPlan-v1.1 | 1B | 4.5M | 60s | 78.00 | 80.91 | 66.38 | 56.85 | 76.30 | 40.35 | 86.80 | 53.11 | 27.17 |
| OpenSora-v1.2 | 1B | 32M | 55s | 79.76 | 81.35 | 73.39 | 56.85 | 82.22 | 51.83 | 91.20 | 68.56 | 42.44 |
| CogVideoX | 2B | 35M | 90s | 80.91 | 82.18 | 75.83 | 60.82 | 83.37 | 62.63 | 98.00 | 69.90 | 51.14 |
| *Autoregressive models* | | | | | | | | | | | | |
| CogVideo | 9B | 5.4M | - | 67.01 | 72.06 | 46.83 | 38.18 | 73.4 | 18.11 | 78.20 | 18.24 | 28.24 |
| Emu3 | 8B | - | - | 80.96 | 84.09 | 68.43 | 59.64 | 86.17 | 44.64 | 77.71 | 68.73 | 37.11 |
| NOVA | 0.6B | 20M | 12s | 78.48 | 78.96 | 76.57 | 54.52 | 91.36 | 73.46 | 91.20 | 66.37 | 50.16 |
| + Rewriter | 0.6B | 20M | 12s | 80.12 | 80.39 | 79.05 | 59.42 | 92.00 | 77.52 | 95.20 | 77.52 | 54.06 |

**NOVA rivals diffusion text-to-video models and significantly suppresses the AR counterpart.** We emphasize that the current version of our NOVA is designed to generate videos at 33 frames and can extend video length through the *pre-filling* of recently generated frames. We perform a quantitative analysis comparing NOVA against open-source and proprietary text-to-video models. As shown in Table 3, despite its significantly smaller size (0.6B vs. 9B), NOVA remarkably outperforms CogVideo (Hong et al. (2022)) across a variety of text-to-video evaluation metrics. It also matches the latest SOTA model Emu3's (Wang et al. (2024)) performance (80.12 vs. 80.96) with a significantly smaller size (0.6B vs. 8B). Additionally, we compared NOVA with state-of-the-art diffusion models. This includes both closed-source models such as Gen-2 (Runway (2023)), Kling (Kuaishou (2024)), and Gen-3 (Runway (2024)), as well as open-source alternatives like LaVie (Wang et al. (2023)), Show-1 (Zhang et al. (2023a)), AnimateDiff-v2 (Guo et al. (2024)), VideoCrafter-v2.0 (Chen et al. (2024a)), T2V-Turbo ( Li et al. (2024b)), OpenSora-v1.1 (Zheng et al. (2024)), OpenSoraPlan-v1.1/v1.2 (Lin et al. (2024)), and CogVideoX (Yang et al. (2024)). The results underscore the effectiveness of text-to-image pre-training within our generalized causal process. Notably, we have narrowed the gap between autoregressive and diffusion methods in modeling large-scale video-text pairs, enhancing both the quality and instruction-following capabilities of video generation. Moreover, NOVA demonstrates a substantial speed advantage over existing models in terms of inference latency.
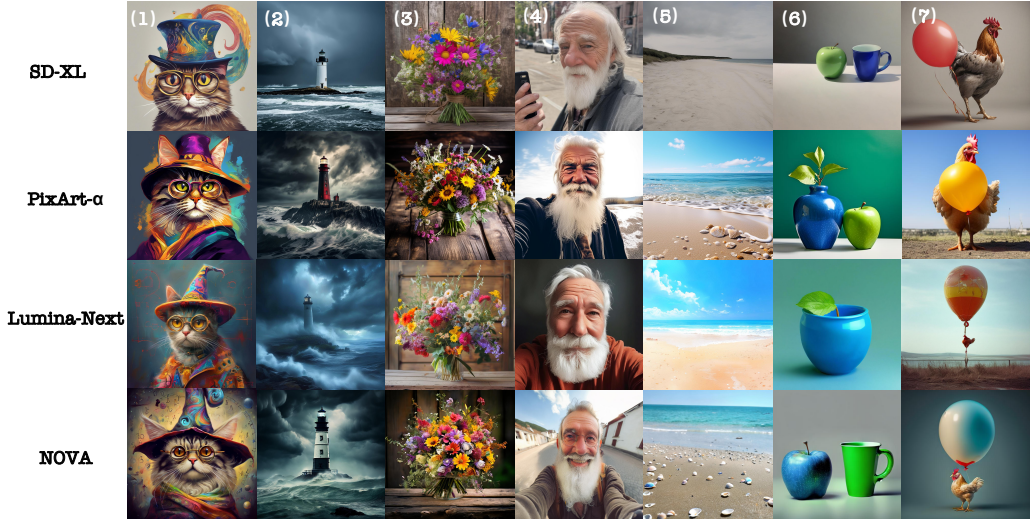
Figure 4: **Text-to-image generation.** Text prompts from left to right: (1) "A digital artwork of a cat styled in a whimsical fashion...", (2) "A solitary lighthouse standing tall against a backdrop of stormy seas and dark, rolling clouds", (3) "A vibrant bouquet of wildflowers on a rustic wooden table", (4) "A selfie of an old man with a white beard", (5) "A serene, expansive beach with no people", (6) "A blue apple and a green cup." and (7) "A chicken on the bottom of a balloon."
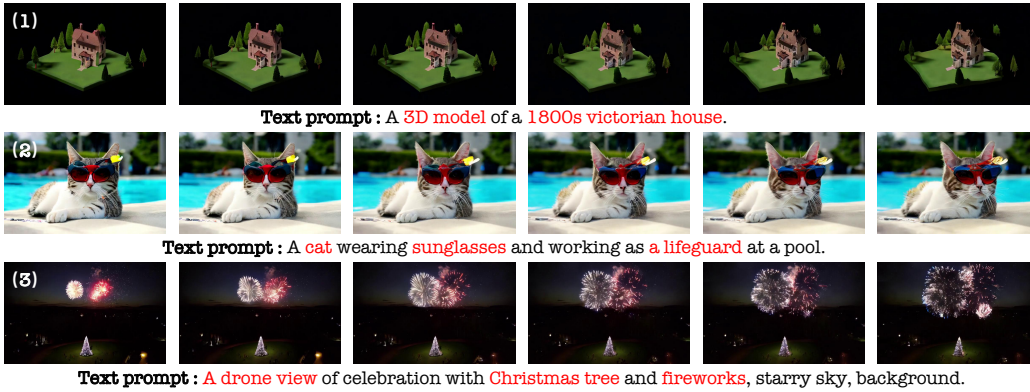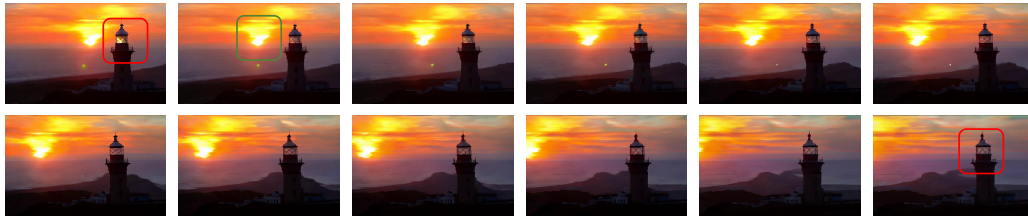


Figure 5: **Text-to-video generation.** We highlight the keywords in red color. NOVA follows the text prompts and vividly captures the motion of subjects (i.e., 3D model, cat and fireworks).

## 4.3 QUALITATIVE RESULTS

**High-fidelity image and high-fluency video.** We present a qualitative comparison of current leading image generation methods in Figure 4. NOVA demonstrates strong visual quality and fidelity across a range of prompt styles, and excels in color attribute binding and spatial object relationships. We present text-to-video visualizations in Figure 5, which highlight NOVA's ability to capture multi-view perspectives, smooth object motion, and stable scene transitions based on the provided text prompts.
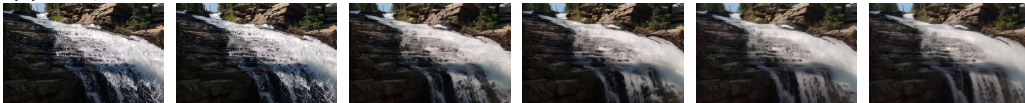
**Zero-shot generalization on video extrapolation.** By pre-filling generated frames, NOVA can produce videos that surpass the training length. For example, by shifting both the text and BOV embeddings, we generate 5-second videos that are up to twice the original length, as shown in Figure 6. We observed that during video extrapolation, NOVA consistently preserves temporal consistency of subject across frames. For instance, when the prompt describes a dome top and a lantern room, the model accurately represents the lighting within the house and captures the transition of a sunset. This further underscores the advantages of causal modeling in long-context video generation tasks.
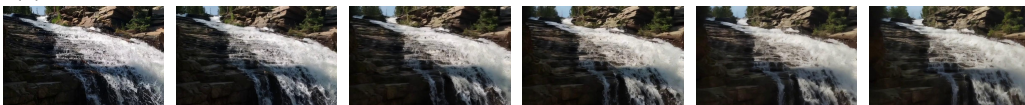
A lighthouse ... a vibrant sunset sky... to a warm palette of oranges and yellows ... has a classic design with a domed top and a lantern room where the light would be housed, ...to intensify.

Figure 6: **Zero-shot video extrapolation.** We highlight the subjects in red and green respectively. The top images are generated, while the bottom images are extrapolated.
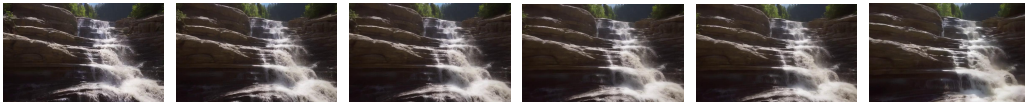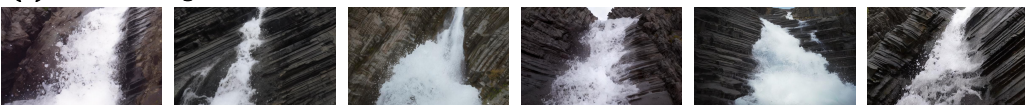
**(1) Image-To-Video Without Text**



**(2) Image-To-Video With Text**



**(3) Text-To-Video**



**(4) Text-To-Image**



**Text :** A cascade of water rushes down a rocky incline, frothing and churning as it descends, is surrounded by rugged, layered rock formations.

Figure 7: **Zero-shot generalization on multiple contexts.** It is evident that NOVA successfully maintains temporal consistency in objects, both with and without text. Such as ensuring "water continues to flow smoothly." This highlights NOVA's capability for zero-shot multitasking.

**Zero-shot generalization on multiple contexts.** By pre-filling the reference image, NOVA can generate videos from images, either with or without accompanying text. In Figure 7, we provide a qualitative example. We show that NOVA can simulate realistic motions without text prompts. Moreover, when text is included, perspective movements appear more natural. This indicates that NOVA is able to capture the fundamental physics, such as interaction forces and fluid dynamics.

## 4.4 ABLATION STUDY

**Effectiveness of temporal autoregressive modeling.** To highlight the advantages of temporal autoregressive modeling, we have facilitated spatial autoregressive to finish video generation task. Specifically, we modify the attention mask of the temporal layer to bidirectional attention, and randomly predict the entire video sequence using set by set prediction. We observe less subject movement in videos under the same training iterations (Figure 8). Additionally, in zero-shot generalization across various contexts or video extrapolation, the network output exhibited more artifacts and temporal inconsistencies. Furthermore, this approach is not compatible with kv-cache acceleration during inference, leading to a linear increase in latency with the number of video frames. This further demonstrates the superiority of causal modeling over multitask approaches for video generation.
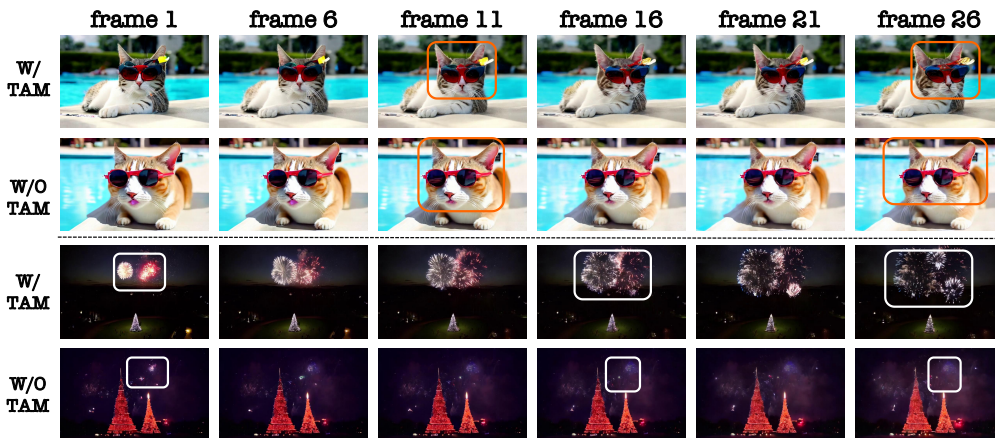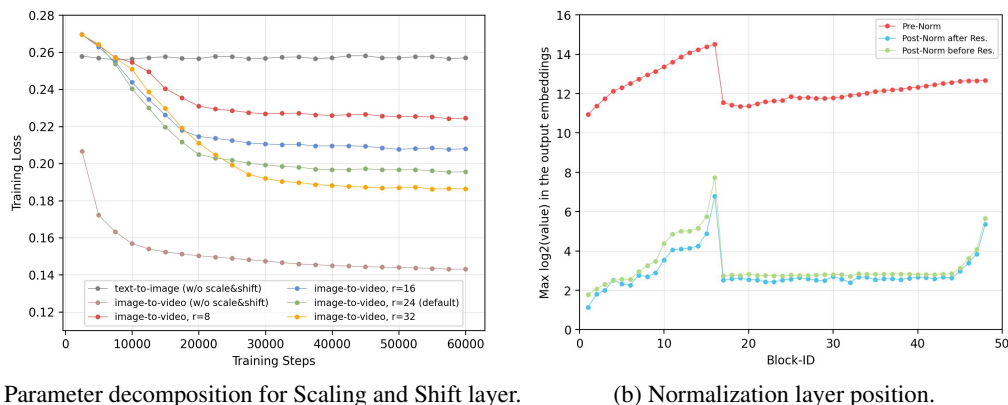
9

Figure 8: **Temporal autoregressive modeling (TAM) for video generation.** We highlight the subtle changes in frames generated from the same prompt. Compared to spatial-only autoregressive method, the inclusion of TAM enables NOVA to more accurately capture the dynamics of subject movement.



(a) Parameter decomposition for Scaling and Shift layer.    (b) Normalization layer position.

Figure 9: **Ablation studies on NOVA's architecture components.** We carefully examine the two key stability factors in large-scale video generation training, as illustrated in (a) and (b).

**Effectiveness of Scaling and Shift Layer.** To capture cross-frame motion changes, we employ a simple yet effective scaling and shifting layer to explicitly model the relative distribution from the BOV-attended feature space. In Figure 9(a), we demonstrate that this approach significantly reduces the drift between text-to-image and image-to-video generation losses. As we gradually decrease the inner rank of the MLP, the training difficulty increases, leading to a more comprehensive and robust learning process for the network. However, extremely low rank values pose challenges for motion modeling, as they significantly limit the layer's representation capability (Figure 10). The rank is set to 24 by default in all text-to-video experiments, resulting in more accurate motion predictions.

**Effectiveness of Post-Norm Layer.** Training large-scale image and video generation models (Ding et al. (2021); ChameleonTeam (2024)) from scratch often poses significant challenges with mixed precision, which is also observed in other visual recognition methods (Liu et al. (2022)). As shown in Figure 9(b), the training process with pre-normalization (Dosovitskiy et al. (2021)) suffers from numerical overflow and variance instability. We attempted various regularization techniques on the residual branch, such as stochastic depth (Huang et al. (2016)) and residual dropout (Vaswani et al. (2017)), but found them to be less effective. Inspired by (Liu et al. (2022)), we introduce post-normalization and empirically discover that it can effectively mitigate the residual accumulation of output embeddings compared to pre-normalization, resulting in a more stable training process.

Figure 10: **Visualization of decomposition ranks in the Scaling and Shift layer.** The first row displays the results of the first frame, while the second row presents the results of the last frame.

## 5 CONCLUSION

In this paper, we present **NOVA**, a novel autoregressive model designed for both text-to-image and text-to-video generation. **NOVA** delivers exceptional image quality and video fluency while significantly minimizing training and inference overhead. Our key designs include temporal frame-by-frame prediction, spatial set-by-set generation, and continuous-space autoregressive modeling across various contexts. Extensive experiments demonstrate that **NOVA** achieves near-commercial quality in image generation, alongside promising fidelity and fluency in video generation. **NOVA** paves the way for next-generation video generation and world models. It offers valuable insights and possibilities for real-time and infinite video generation, going beyond Sora-like video diffusion models. As a first step, we will continue scalable experiments with larger models and data scaling to explore NOVA's limits in future work.

## REFERENCES

Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.

Jason Baldridge, Jakob Bauer, Mukul Bhutani, Nicole Brichtova, Andrew Bunner, Kelvin Chan, Yichang Chen, Sander Dieleman, Yuqing Du, Zach Eaton-Rosen, et al. Imagen 3. *arXiv preprint arXiv:2408.07009*, 2024.

James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science. https://cdn. openai. com/papers/dall-e-3. pdf*, 2(3):8, 2023a.

James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science. https://cdn. openai. com/papers/dall-e-3. pdf*, 2(3):8, 2023b.

Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.

G. Bradski. The OpenCV Library. *Dr. Dobb's Journal of Software Tools*, 2000.

Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, et al. Video generation models as world simulators, 2024. URL https://openai.com/research/video-generation-modelsas-world-simulators.

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. In *Advances in Neural Information Processing Systems*, 2020.

Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Saehoon Kim. Coyo-700m: Image-text pair dataset. https://github.com/kakaobrain/coyo-dataset, 2022.

ChameleonTeam. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024.

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*, pp. 11315–11325, 2022.

Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*, 2023.

Haoxin Chen, Yong Zhang, Xiaodong Cun, Menghan Xia, Xintao Wang, Chao Weng, and Ying Shan. Videocrafter2: Overcoming data limitations for high-quality video diffusion models, 2024a.

Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart-$\alpha$: Fast training of diffusion transformer for photorealistic text-to-image synthesis, 2023.

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. In *European Conference on Computer Vision*, pp. 74–91. Springer, 2024b.

Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Ekaterina Deyneka, Hsiang-wei Chao, Byung Eun Jeon, Yuwei Fang, Hsin-Ying Lee, Jian Ren, Ming-Hsuan Yang, et al. Panda-70m: Captioning 70m videos with multiple cross-modality teachers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13320–13331, 2024c.

Haiwen Diao, Yufeng Cui, Xiaotong Li, Yueze Wang, Huchuan Lu, and Xinlong Wang. Unveiling encoder-free vision-language models. *arXiv preprint arXiv:2406.11832*, 2024.

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

Ming Ding, Wendi Zheng, Wenyi Hong, and Jie Tang. Cogview2: Faster and better text-to-image generation via hierarchical transformers. *Advances in Neural Information Processing Systems*, 35: 16890–16902, 2022.

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 *ICLR*, 2021.

Stefan Elfwing, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. *Neural networks*, 107:3–11, 2018.

Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 12873–12883, 2021.

Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germanidis. Structure and content-guided video synthesis with diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7346–7356, 2023.

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*, 2024a.

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*, 2024b.

Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruba Ghosh, Jieyu Zhang, et al. Datacomp: In search of the next generation of multimodal datasets. *Advances in Neural Information Processing Systems*, 2024.

Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36, 2024.

P Goyal. Accurate, large minibatch sg d: training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.

Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023.

Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. *International Conference on Learning Representations*, 2024.

Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross B. Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, pp. 15979–15988, 2022.

Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.

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

Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022.

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

Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip diffusion models with llm for enhanced semantic alignment. *arXiv preprint arXiv:2403.05135*, 2024.

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*, pp. 646–661. Springer, 2016.

Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2i-compbench: A comprehensive benchmark for open-world compositional text-to-image generation. *Advances in Neural Information Processing Systems*, 2023.

Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE international conference on computer vision*, pp. 1501–1510, 2017a.

Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE international conference on computer vision*, pp. 1501–1510, 2017b.

Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21807–21818, 2024.

Mojan Javaheripi, Sébastien Bubeck, Marah Abdin, Jyoti Aneja, Sebastien Bubeck, Caio César Teodoro Mendes, Weizhu Chen, Allie Del Giorno, Ronen Eldan, Sivakanth Gopi, et al. Phi-2: The surprising power of small language models. *Microsoft Research Blog*, 2023.

Nal Kalchbrenner, Aäron Oord, Karen Simonyan, Ivo Danihelka, Oriol Vinyals, Alex Graves, and Koray Kavukcuoglu. Video pixel networks. In *International Conference on Machine Learning*, pp. 1771–1779. PMLR, 2017.

Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4401–4410, 2019.

Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, Jonathan Huang, Rachel Hornung, Hartwig Adam, Hassan Akbari, Yair Alon, Vighnesh Birodkar, et al. Videopoet: A large language model for zero-shot video generation. *arXiv preprint arXiv:2312.14125*, 2023.

Kuaishou. Kling ai, 2024. URL https://klingai.com/.

Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *ArXiv e-prints*, pp. arXiv–1607, 2016.

Daiqing Li, Aleks Kamko, Ehsan Akhgari, Ali Sabet, Linmiao Xu, and Suhail Doshi. Playground v2. 5: Three insights towards enhancing aesthetic quality in text-to-image generation. *arXiv preprint arXiv:2402.17245*, 2024a.

Jiachen Li, Weixi Feng, Tsu-Jui Fu, Xinyi Wang, Sugato Basu, Wenhu Chen, and William Yang Wang. T2v-turbo: Breaking the quality bottleneck of video consistency model with mixed reward feedback. *arXiv preprint arXiv:2405.18750*, 2024b.

Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image generation without vector quantization. *arXiv preprint arXiv:2406.11838*, 2024c.

Zhimin Li, Jianwei Zhang, Qin Lin, Jiangfeng Xiong, Yanxin Long, Xinchi Deng, Yingfang Zhang, Xingchao Liu, Minbin Huang, Zedong Xiao, et al. Hunyuan-dit: A powerful multi-resolution diffusion transformer with fine-grained chinese understanding. *arXiv preprint arXiv:2405.08748*, 2024d.

Jun Hao Liew, Hanshu Yan, Jianfeng Zhang, Zhongcong Xu, and Jiashi Feng. Magicedit: High-fidelity and temporally coherent video editing. *arXiv preprint arXiv:2308.14749*, 2023.

Bin Lin, Yunyang Ge, Xinhua Cheng, Zongjian Li, Bin Zhu, Shaodong Wang, Xianyi He, Yang Ye, Shenghai Yuan, Liuhan Chen, et al. Open-sora plan: Open-source large video generation model, 2024.

Bingchen Liu, Ehsan Akhgari, Alexander Visheratin, Aleks Kamko, Linmiao Xu, Shivam Shrirao, Chase Lambert, Joao Souza, Suhail Doshi, and Daiqing Li. Playground v3: Improving text-to-image alignment with deep-fusion large language models. *arXiv preprint arXiv:2409.10695*, 2024.

Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 12009–12019, 2022.

Ilya Loshchilov, Frank Hutter, et al. Fixing weight decay regularization in adam. *arXiv preprint arXiv:1711.05101*, 2017.

Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021.

Charlie Nash, Joao Carreira, Jacob Walker, Iain Barr, Andrew Jaegle, Mateusz Malinowski, and Peter Battaglia. Transframer: Arbitrary frame prediction with generative models. *arXiv preprint arXiv:2203.09494*, 2022.

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.

Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International conference on machine learning*, pp. 8162–8171. PMLR, 2021.

William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.

Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.

PexelsTeam. Pexels, royalty-free stock footage website, 2014.

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. *arXiv preprint arXiv:2307.01952*, 2023.

Alec Radford. Improving language understanding by generative pre-training, 2018.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

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*, pp. 8821–8831. Pmlr, 2021.

Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.

Scott Reed, Aäron Oord, Nal Kalchbrenner, Sergio Gómez Colmenarejo, Ziyu Wang, Yutian Chen, Dan Belov, and Nando Freitas. Parallel multiscale autoregressive density estimation. In *International conference on machine learning*, pp. 2912–2921. PMLR, 2017.

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*, pp. 10684–10695, 2022a.

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*, pp. 10684–10695, 2022b.

Runway. Gen-2: Generate novel videos with text, images or video clips, 2023. URL https://runwayml.com/research/gen-2.

Runway. Gen-3 alpha: A new frontier for video generation, 2024. URL https://runwayml.com/research/introducing-gen-3-alpha.

Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *NeurIPS*, 35:25278–25294, 2022.

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.

Keqiang Sun, Junting Pan, Yuying Ge, Hao Li, Haodong Duan, Xiaoshi Wu, Renrui Zhang, Aojun Zhou, Zipeng Qin, Yi Wang, et al. Journeydb: A benchmark for generative image understanding. *Advances in Neural Information Processing Systems*, 2024a.

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*, 2024b.

Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiying Yu, Zhengxiong Luo, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative multimodal models are in-context learners. *arXiv: 2312.13286*, 2023.

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.

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.

UnsplashTeam. Unsplash dataset, 2020.

Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, pp. 5998–6008, 2017.

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 descriptions. In *International Conference on Learning Representations*, 2022.

Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, Quan Sun, Yufeng Cui, Jinsheng Wang, Fan Zhang, Yueze Wang, Zhen Li, Qiying Yu, et al. Emu3: Next-token prediction is all you need. *arXiv preprint arXiv:2409.18869*, 2024.

Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yinan He, Jiashuo Yu, Peiqing Yang, et al. Lavie: High-quality video generation with cascaded latent diffusion models. *arXiv preprint arXiv:2309.15103*, 2023.

Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. Videogpt: Video generation using vq-vae and transformers. *arXiv preprint arXiv:2104.10157*, 2021.

Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024.

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. *arXiv preprint arXiv:2206.10789*, 2(3):5, 2022.

Lijun Yu, Yong Cheng, Kihyuk Sohn, José Lezama, Han Zhang, Huiwen Chang, Alexander G Hauptmann, Ming-Hsuan Yang, Yuan Hao, Irfan Essa, et al. Magvit: Masked generative video transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10459–10469, 2023.

David Junhao Zhang, Jay Zhangjie Wu, Jia-Wei Liu, Rui Zhao, Lingmin Ran, Yuchao Gu, Difei Gao, and Mike Zheng Shou. Show-1: Marrying pixel and latent diffusion models for text-to-video generation. *arXiv preprint arXiv:2309.15818*, 2023a.

Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3836–3847, 2023b.

Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 586–595, 2018.

Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all, March 2024. URL https://github.com/hpcaitech/Open-Sora.

Le Zhuo, Ruoyi Du, Han Xiao, Yangguang Li, Dongyang Liu, Rongjie Huang, Wenze Liu, Lirui Zhao, Fu-Yun Wang, Zhanyu Ma, et al. Lumina-next: Making lumina-t2x stronger and faster with next-dit. *arXiv preprint arXiv:2406.18583*, 2024.

## APPENDIX

We strictly publish our code and pretrained models to improve interpretability and assure reproducibility. Here, more implementation details and ablation experiments are organized as follows:

- Architecture details of Scaling and Shift layer (Sec. A)
- Normalization configurations (Sec. B)
- Video extrapolation evaluations (Sec. C)
- Inference time analysis (Sec. D)
- Ablations on the impact of temporal autoregressive modeling (Sec. E)
- Comprehensive DPG-Bench evaluation results (Sec. F)
- More text-to-image visualizations (Sec. G)
- More text-to-video visualizations (Sec. H)

## A    ARCHITECTURE DETAILS OF SCALING AND SHIFT LAYER

The Scaling and Shift Layer is implemented as an adaptive normalization layer, adopting the design initially proposed by FiLM (Perez et al. (2018)) and AdaIN (Huang & Belongie (2017b)). While many previous methods have primarily utilized adaptive normalization for controllable image generation, such as in StyleGAN (Karras et al. (2019)), or for conditional modeling within Diffusion Transformers, like DiT (Peebles & Xie (2023)), NOVA innovatively applies this technique to manage the cumulative inference errors in autoregressive video generation. We employ a two-layer MLP to optimize low-rank decomposition for motion changes, as shown in the Figure 11. Specifically, we refer `AdaLayerNorm` and decompose the motion changes into mean and variance parameters, which are further used to apply the affine transformation on BOV embeddings.
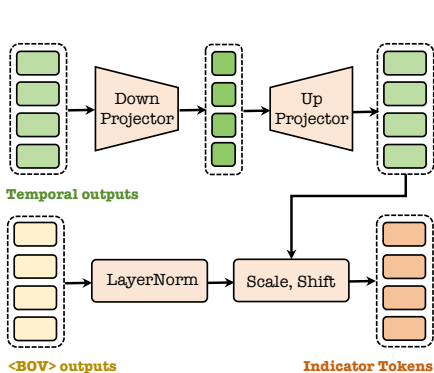


Figure 11: **Scaling and Shift layer.** We reformulate cross-frame motion changes by learning relative distribution variations within a unified space based on BOV tokens, rather than directly modeling the unreferenced distribution of the current frame.
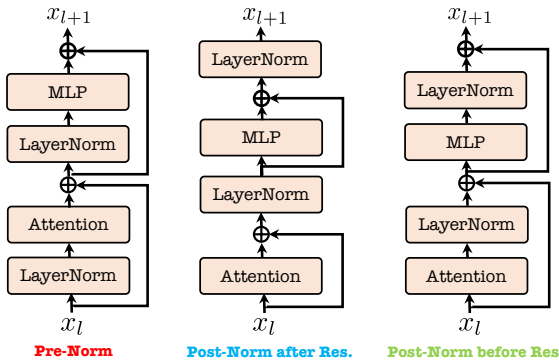
Figure 12: **Three normalization architectures.** We summarize various configurations including the pre-normalization layer (left), the post-normalization layer after residual addition (middle), and the post-normalization layer before residual addition (right). Here Post-Norm before Res is our standard design.

## B    NORMALIZATION CONFIGURATIONS

NOVA employs an improved normalization configuration that can effectively control the numerical boundaries of the output embeddings of each Transformer block while also maintaining the identity transformation of residual connections. We illustrate the three common normalization configurations in Figure 12, and NOVA uses the *post-normalization before residual addition* by default.
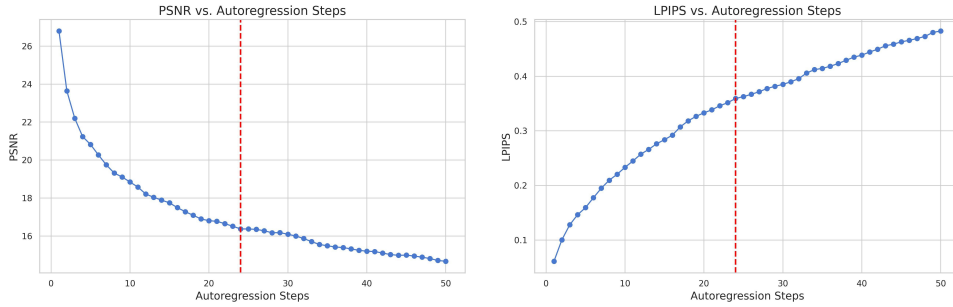
Figure 13: **PSNR and LPIPS metrics over 50 autoregressive steps in video extrapolation.** Due to the 4× downsampling rate of VAE in temporal scale, each autoregressive step generates four frames. The vertical red line marks the point where the extrapolation reaches 3× training length.



Figure 14: **Visualization of video extrapolation.** Although the metrics indicate a decline, the generated frames still closely resemble the original video in content and overall image quality. Visualization suggests that the model can extrapolate up to 3× training length.

## C   VIDEO EXTRAPOLATION EVALUATIONS

Video extrapolation represents a significant challenge, being an out-of-domain generalization issue. To assess our model's performance, we curated a test set comprising 200 videos. For each video, the task involved generating subsequent frames from the initial frame and a textual prompt, effectively converting an image and text into a video sequence. We utilized LPIPS (Zhang et al. (2018)) and PSNR metrics to evaluate the video extrapolation capabilities of our model.

During the extrapolation process, it was observed that the generated frames started to deviate from the ground truth after a few iterations. This is mainly due to the difficulty in accurately capturing video dynamics, causing minor discrepancies to accumulate. As a result, per-frame PSNR values decrease, while LPIPS scores increase over time (Figure 13). Nevertheless, the generated frames exhibit a high degree of similarity to the original video in terms of both content and image quality in Figure 14. This highlights the robustness of our temporal autoregressive approach in video extrapolation.

## D   INFERENCE TIME ANALYSIS

We report inference times on a single NVIDIA A100 GPU (40GB) with a batch size of 24 in Table 4. In each video, the temporal layers require only 0.03 seconds, compared to 11.97 seconds for the spatial layers, highlighting the exceptional efficiency of the temporal layers. While NOVA is already efficient in text-to-video generation, there is potential for further acceleration in the spatial layers.

Table 4: **Inference time analysis for different layers.**

| Resolution | Temporal Layers Time | Spatial Layers Time | Total Time |
|---|---|---|---|
| 29×768×480 | 0.03s | 11.97s | 12s |

# E  ABLATIONS ON THE IMPACT OF TEMPORAL AUTOREGRESSIVE MODELING

Under the same settings, we evaluate VBench results in Table 5 with and without using TAM (Temporal Autoregressive Modeling) to highlight its significance. Our findings are summarized as follows: **(1) Efficient Motion Modeling:** We observed that the total score was marginally lower without TAM compared to NOVA (75.38 vs. 75.84), especially in the dynamic degree metric, which showed a more pronounced decline (11.38 vs. 23.27). We hypothesize that while bidirectional attention enhances model capacity, it requires more extensive data and longer training times to capture subtle motion changes compared to causal models. **(2) Efficient Video Inference:** Thanks to the kv-cache technology and frame-by-frame autoregressive processing, NOVA's inference time is much faster compared to methods without TAM, with a greater speed advantage for longer videos.

Table 5: **Performance comparison on temporal autoregressive modeling.**

| Model | Total Score | Dynamic Degree | Infer Time |
|---|---|---|---|
| NOVA | 75.84 | 23.27 | 12s |
| NOVA (w/o TAM) | 75.38 | 11.38 | 39s |

# F  COMPREHENSIVE DPG-BENCH EVALUATION RESULTS

We provide detailed DPG-Bench scores in Table 6. While NOVA outperforms most models of comparable size and matches the overall score of state-of-the-art models, we observe that increasing the model scale results in marginal improvements and does not boost the text rendering performance. This limitation may be attributed to our reliance on extensive web datasets, such as LAION and DataComp. In future work, we plan to focus on improving the quality of text-to-image data.

Table 6: **Comparison with state-of-the-art models on DPG-Bench.**

| Model | Overall | Global | Entity | Attribute | Relation | Other |
|---|---|---|---|---|---|---|
| *Diffusion models* | | | | | | |
| SD v1.5 (Rombach et al. (2022b)) | 63.18 | 74.63 | 74.23 | 75.39 | 73.49 | 67.81 |
| PixArt-$\alpha$ (Chen et al. (2023)) | 71.11 | 74.97 | 79.32 | 78.60 | 82.57 | 76.96 |
| PixArt-$\sigma$ (Chen et al. (2024b)) | 80.54 | 86.89 | 82.89 | 88.94 | 86.59 | 87.68 |
| Lumina-Next (Zhuo et al. (2024)) | 74.63 | 82.82 | 88.65 | 86.44 | 80.53 | 81.82 |
| SDXL (Podell et al. (2023)) | 74.65 | 83.27 | 82.43 | 80.91 | 86.76 | 80.41 |
| Playground v2.5 (Li et al. (2024a)) | 75.47 | 83.06 | 82.59 | 81.20 | 84.08 | 83.50 |
| Hunyuan-DiT (Li et al. (2024d)) | 78.87 | 84.59 | 80.59 | 88.01 | 74.36 | 86.41 |
| DALL-E3 (Betker et al. (2023a)) | 83.50 | 90.97 | 89.61 | 88.39 | 90.58 | 89.83 |
| SD3 (Esser et al. (2024b)) | 84.08 | 87.90 | 91.01 | 88.83 | 80.70 | 88.68 |
| Playground v3 (Liu et al. (2024)) | 87.04 | 91.94 | 85.71 | 90.90 | 90.00 | 92.72 |
| *Autoregressive models* | | | | | | |
| Emu3-DPO (Wang et al. (2024)) | 81.60 | 87.54 | 87.17 | 86.33 | 90.61 | 89.75 |
| NOVA (0.3B) | 80.60 | 85.41 | 86.97 | 85.16 | 92.05 | 71.20 |
| NOVA (0.6B) | 82.25 | 87.65 | 87.65 | 85.62 | 90.90 | 74.80 |
| NOVA (1.4B) | 83.01 | 86.32 | 88.69 | 86.35 | 91.94 | 74.80 |

We present more text-to-image samples in the Figure 15. NOVA can generate images with a maximum resolution of 1024×1024. Our model excels in the domain of text-to-image generation, producing a vast array of high-quality images that accurately reflect the textual descriptions provided. This capability not only spans a wide range of subjects, from realistic landscapes and portraits to imaginative and abstract concepts, but also maintains a high level of detail and aesthetic quality.
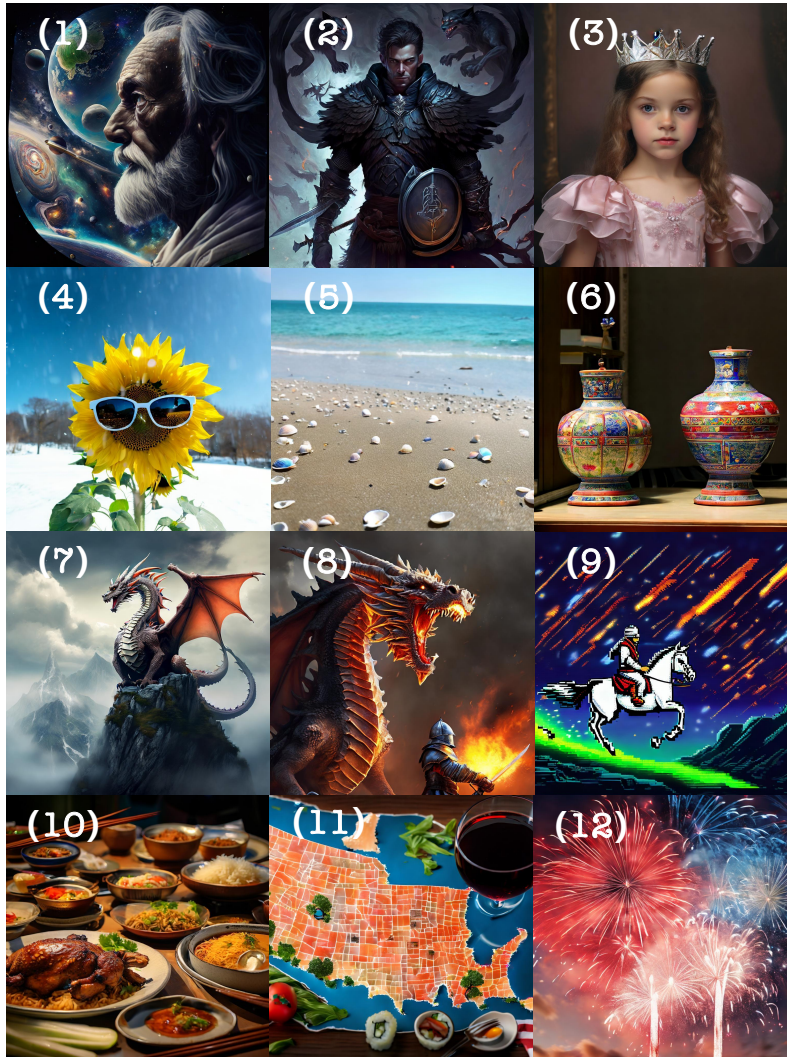


Figure 15: **More text-to-image visualizations.** Text prompts are as follows: (1) "In the foreground is the detailed, head-and-shoulders portrait of an elderly man with a long white beard...", (2) "a digital artwork of a fantasy warrior character. The character is male, depicted from the waist up, and appears to have a stern or serious facial expression...", (3) "a young girl wearing a tiara and frilly dress", (4) "A sunflower in sunglasses dances in the snow", (5) "A beach with no people", (6) "Two Ming vases on the table, the larger one is more colorful than the other", (7) "A dragon perched majestically on a craggy, smoke-wreathed mountain", (8) "a dragon breathing fire onto a knight", (9) "a pixel art style graphic with vibrant colors. It features a single rider on a horse, both depicted in mid-gallop to the left side of the frame...", (10) "A table full of food. There is a plate of chicken rice, a bowl of bak chor mee, and a bowl of laksa", (11) "A map of the United States made out sushi. It is on a table next to a glass of red wine" and (12) "beautiful fireworks in the sky with red, white and blue".

# H   MORE TEXT-TO-VIDEO VISUALIZATIONS

We present more text-to-video samples generated by NOVA in the Figure 16. NOVA can generate videos with a resolution of 33×768×480. Our model stands out in the field of text-to-video generation, capable of producing a substantial number of high-quality videos that vividly bring textual descriptions to life. From detailed storylines and character animations to realistic environmental settings and action scenes, our model demonstrates exceptional proficiency in generating content.



**(1)** Text prompt : A 3D model of a 1800s victorian house.

**(2)** Text prompt : A dog drinking water.

**(3)** Text prompt : A drone view of celebration with Christmas tree and fireworks, starry sky, background.

**(4)** Text prompt : A space shuttle launching into orbit, with flames and smoke billowing out from the engines.

**(5)** Text prompt : A 3D model of a 1800s victorian house.

**(6)** Text prompt : Extreme close-up of chicken and green pepper kebabs grilling on a barbeque with flames. Shallow focus and light smoke. vivid colours..
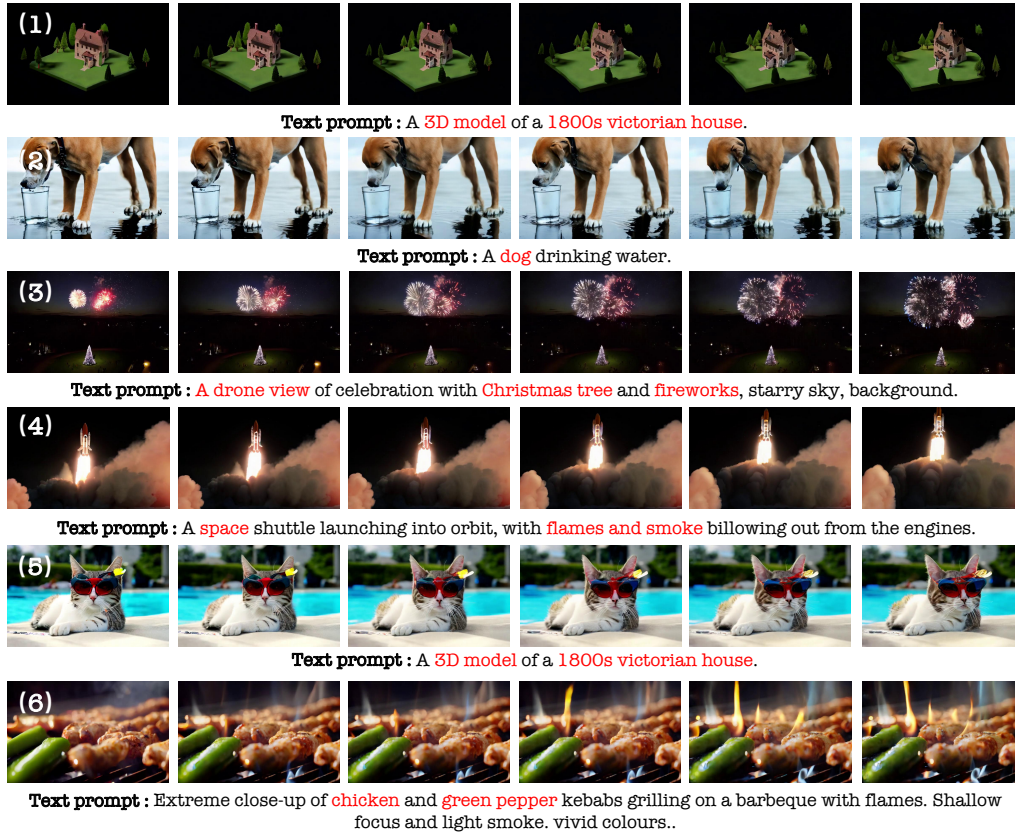
Figure 16: More text-to-video visualizations. Best viewed with zoom for enhanced detail.