# From Reusing to Forecasting: Accelerating Diffusion Models with TaylorSeers

Jiacheng Liu[1,2*], Chang Zou[1,3*], Yuanhuiyi Lyu[4], Junjie Chen[1], Linfeng Zhang[1†]

[1]Shanghai Jiao Tong University   [2]Shandong University
[3]University of Electronic Science and Technology of China
[4]The Hong Kong University of Science and Technology

## Abstract

*Diffusion Transformers (DiT) have revolutionized high-fidelity image and video synthesis, yet their computational demands remain prohibitive for real-time applications. To solve this problem, feature caching has been proposed to accelerate diffusion models by caching the features in the previous timesteps and then reusing them in the following timesteps. However, at timesteps with significant intervals, the feature similarity in diffusion models decreases substantially, leading to a pronounced increase in errors introduced by feature caching, significantly harming the generation quality. To solve this problem, we propose TaylorSeer, which firstly shows that features of diffusion models at future timesteps can be predicted based on their values at previous timesteps. Based on the fact that features change slowly and continuously across timesteps, TaylorSeer employs a differential method to approximate the higher-order derivatives of features and predict features in future timesteps with Taylor series expansion. Extensive experiments demonstrate its significant effectiveness in both image and video synthesis, especially in high acceleration ratios. For instance, it achieves an almost lossless acceleration of $4.99\times$ on FLUX and $5.00\times$ on HunyuanVideo without additional training. On DiT, it achieves $3.41$ lower FID compared with previous SOTA at $4.53\times$ acceleration. Our codes have been released in Github:*
***https://github.com/Shenyi-Z/TaylorSeer***

## 1. Introduction

Diffusion Models (DMs) [11] have made significant strides in generative artificial intelligence, achieving remarkable results in tasks such as image generation and video synthesis [1, 31]. The introduction of Diffusion Transformers (DiT) [28] has further advanced the quality of visual genera-
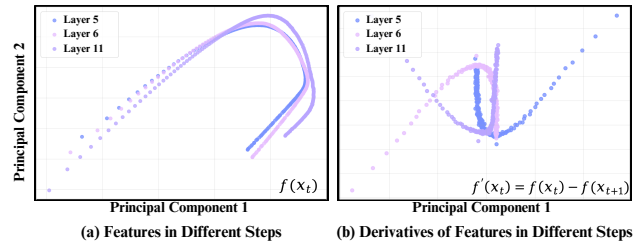


Figure 1. **PCA projections of features and their derivatives in diffusion models.** Both the features (a) and derivatives of features (b) in diffusion models at different timesteps form stable trajectories, indicating it possible to predict features of diffusion models at future timesteps based on features from previous timesteps.

tion. However, these improvements come with a substantial increase in computational demands, limiting the practical use of diffusion transformers. In response to these growing challenges regarding computational efficiency, various acceleration techniques have been proposed [26, 46, 49, 53]. Recently, based on the observation that diffusion models have highly similar features in the adjacent timesteps, feature caching methods have been proposed to store the features in the previous timesteps and then reuse them in the following timesteps, allowing diffusion models to skip substantial computations in both U-Net-based and transformer-based diffusion models [5, 26, 37] without requirements for additional training. Previous caching methods follow the *"cache-then-reuse"* paradigm and explore it from the perspectives of tokens [37, 53, 54] and residuals [5] while ignoring the its natural limitation: *As the distance between two timesteps increases, their feature similarity decreases exponentially, making reusing features at distant steps significantly harm the generation quality*. This curse locks the possibility of feature caching in high-ratio acceleration. As shown in the PCA results from Figure 1(a), features of diffusion models in non-adjacent timesteps exhibit long distances, verifying this limitation and calling for the development of a new caching paradigm.

To solve this problem, this paper introduces a new paradigm, *"cache-then-forecast"*, to replace the previous *"cache-then-reuse"*. As shown in Figure 1(a), features of

---
*Equal contribution. shenyizou@outlook.com
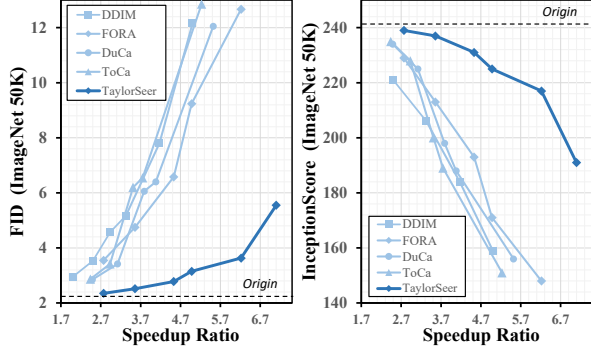†Corresponding author: zhanglinfeng@sjtu.edu.cn

Figure 2. **Comparison between previous caching methods and TaylorSeer.** TaylorSeer shows significantly better performance at high-acceleration ratios.

the diffusion model at different timesteps form a stable trajectory, demonstrating that it is possible to predict the features at the future timesteps based on features from previous timesteps. Furthermore, we investigate the derivatives of features in different timesteps (*i.e.,* the velocity along the trajectory) and present their PCA results in Figure 1(b). Surprisingly, derivatives of features also exhibit similar values at the adjacent timesteps with high stability, indicating that predicting future features is not a complex problem and may be solvable using a non-parametric method.

Building on this observation, we propose *TaylorSeer*, a *"cache-then-forecast"* paradigm that leverages Taylor series expansion to predict the features at future timesteps. Specifically, *TaylorSeer* utilizes multi-step features to approximate derivatives of various orders as features evolve over timesteps using difference methods and applies the Taylor series to predict subsequent features. Unlike previous methods that directly reuse cached features, our approach exploits the continuity of feature changes to predict future features, allowing diffusion models to achieve a training-free and high-ratio acceleration without significant decrements in generation quality. Compared with previous caching methods that suffer from low feature similarity at distant timesteps, *TaylorSeer* is particularly effective at larger intervals between full activations, where higher-order Taylor series approximations excel in long-range feature reuse scenarios. As shown in Figure 2, our method reduces quality loss compared with the previous SOTA by 36×, achieving promising performance in acceleration regimes beyond 6×, where all previous methods fail.

In summary, our contributions are as follows:

- **Cache-then-forecast Paradigm**: We propose a novel paradigm of *"cache-then-forecast"* to replace the previous *"cache-then-resue"*, enabling the prediction of diffusion model features at future timesteps through sequential modeling and overcoming the limitations of caching-based methods in high-acceleration regimes.
- **TaylorSeer**: We introduce *TaylorSeer*, which utilizes

Taylor series expansion to approximate the trajectory of features at different timesteps using higher-order derivatives, surpassing previous methods by a large margin without introducing any training or search costs.
- **State-of-the-Art Performance**: *TaylorSeer* achieves **2.5×**, **4.99×**, and **5.00×** acceleration on DiT, FLUX, and HunyuanVideo for image and video synthesis, respectively, while maintaining high-quality generation and sometimes even providing additional benefits, paving a new path for diffusion model acceleration.

## 2. Related Works

Diffusion models [12, 39] have demonstrated exceptional capabilities in image and video generation. Initial architectures primarily utilized U-Net-based designs [32] which, despite their efficacy, encountered scalability constraints that limited larger model training and practical deployment. The introduction of Diffusion Transformer (DiT) [29] addressed these limitations, subsequently inspiring numerous advancements that adapted the architecture to achieve state-of-the-art performance across diverse domains [3, 4, 45, 51]. Despite these achievements, the inherent sequential sampling process of diffusion models imposes substantial computational demands during inference. Consequently, acceleration techniques have emerged as a critical research focus, broadly categorized as *Sampling Timestep Reduction* and *Denoising Network Acceleration*.

### 2.1. Sampling Timestep Reduction

A fundamental approach to accelerating diffusion models involves *minimizing sampling steps while preserving output quality*. DDIM [40] established a deterministic sampling methodology that maintained generation fidelity with reduced denoising iterations. The DPM-Solver series [24, 25, 50] further advanced this concept through high-order ODE solvers. Alternative strategies include Rectified Flow [22], which constructs direct paths between noise and data distributions, and knowledge distillation techniques [27, 36] that compress multiple denoising operations into fewer steps. Notably, Consistency Models [41] introduced an innovative framework enabling single-step or few-step sampling by directly mapping noisy inputs to clean data, eliminating sequential denoising requirements and significantly enhancing practical applicability.

### 2.2. Denoising Network Acceleration

In addition to reducing the number of sampling timesteps, *optimizing the computational efficiency of the denoising network itself* presents another promising approach for accelerating inference. This can be broadly classified into *Model Compression-based* and *Feature Caching-based* techniques, as detailed below.
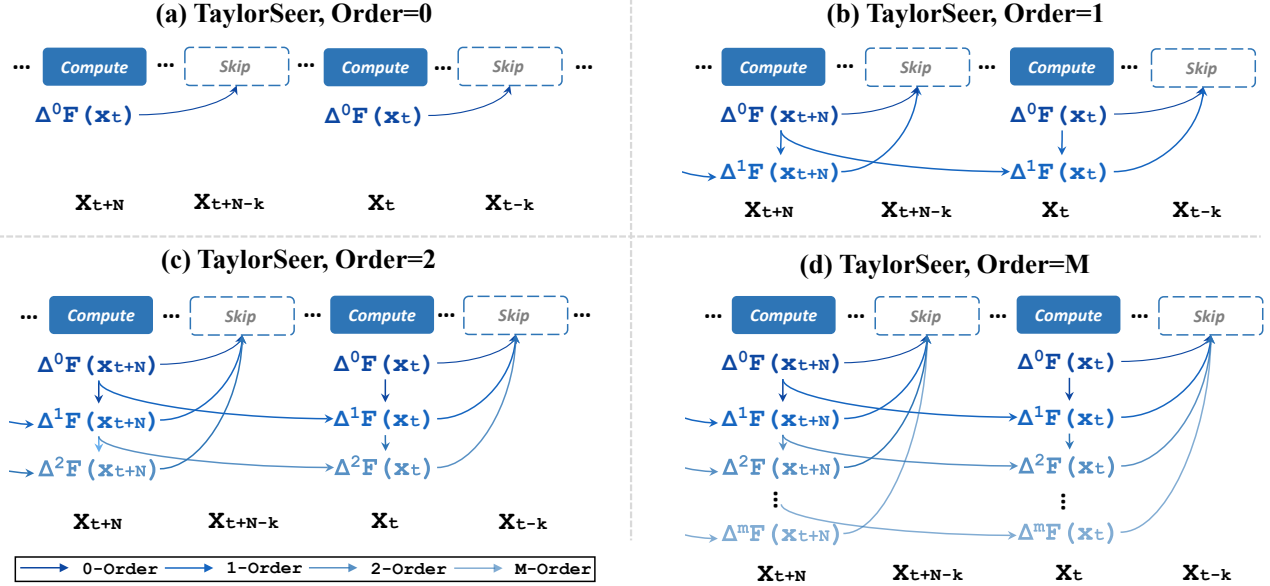
Figure 3. **An overview of TaylorSeer**. (a) **TaylorSeer (Order=0)** *Naïve feature caching*, which directly reuses computed features across timesteps. (b) **TaylorSeer (Order=1)** *Linear prediction*, which estimates feature trajectories using first-order finite differences. (c) **TaylorSeer (Order=2)** extends linear prediction to 2-order finite differences for more accurate modeling of nonlinear feature trajectories. (d) **TaylorSeer (Order=M)** further extends to M-orders for improved accuracy without sacrificing efficiency.

**Model Compression-based Acceleration.** Model compression techniques encompass network pruning [8, 52], quantization [15, 18, 38], knowledge distillation [19], and token reduction [2, 6, 14, 34, 47, 48]. These methods typically require additional training to fine-tune the compressed model, ensuring minimal degradation in generation quality while significantly improving inference speed. Despite their effectiveness, these approaches often involve a trade-off between model size reduction and potential loss of expressive power, making it essential to design efficient compression strategies tailored to diffusion models.

**Feature Caching-based Acceleration.** Feature caching has emerged as particularly advantageous due to its training-free implementation. Initially developed for U-Net architectures [17, 26], caching mechanisms have evolved to address the computational demands of DiT models, which offer superior performance but at increased computational cost. Advanced techniques such as FORA [37] and $\Delta$-DiT [5] leverage attention and MLP representation reuse, while TeaCache [21] implements dynamic timestep-dependent difference estimation to optimize caching decisions. DiTFastAttn [46] reduces redundancies in self-attention computation across spatial, temporal, and conditional dimensions. The ToCa series [53, 54] mitigates information loss through dynamic feature updates, while EOC [30] introduces an error-optimized framework utilizing prior knowledge extraction and adaptive optimization. Recent innovations including UniCP [42], which in-

tegrates dynamic cache window adjustment with pruning, and RAS [23], which implements region-adaptive sampling rates based on model focus, further enhance computational efficiency while preserving generation fidelity.

Despite these advancements, existing caching methods predominantly follow a *"cache-then-reuse"* paradigm, which stores and reuses features based on tokens and residuals. However, as the timestep gap increases, feature similarity decreases exponentially, leading to degradation in generation quality and limiting the potential acceleration gains. In this work, we focus on advancing feature caching techniques by introducing a novel *"cache-then-forecast"* paradigm. Instead of directly reusing stored features, our approach predicts future features based on past ones, leveraging their stable trajectory across timesteps. This predictive strategy enables significantly higher acceleration while maintaining generation fidelity, all without relying on complex additional models.

## 3. Method

### 3.1. Preliminary

**Diffusion Models.** Diffusion models operate through a continuous-time stochastic process with forward and reverse phases for precise noise addition and removal, governed by stochastic differential equations (SDEs). The forward diffusion process is mathematically defined as:

$$dx = \sqrt{\frac{d\sigma^2(t)}{dt}}dw, \tag{1}$$

where $\sigma(t)$ denotes the noise schedule, and $w$ represents the standard Wiener process. The corresponding reverse process restores data by removing noise:

$$dx = \left[-\frac{1}{2}\beta(t)x - \nabla_{x_t} \log p_t(x_t)\right] dt + \sqrt{\beta(t)}d\bar{w}. \quad (2)$$

This formulation ensures that features evolve smoothly over time in a continuous manner. In practical implementations, the process is discretized into timesteps, maintaining its structured trajectory. Specifically, the denoising transition at each step follows a Gaussian distribution:

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}\left(x_{t-1}; \mu_\theta(x_t, t), \beta_t\mathbf{I}\right). \quad (3)$$

Despite the inherent discretization challenges, the underlying feature transformations remain predictable due to the fundamentally smooth nature of the diffusion trajectory.

**Assumption 1.** *The feature representations in diffusion models evolve smoothly over time. Specifically, the underlying feature transformation is a differentiable function with bounded higher-order derivatives, ensuring structured and predictable variation across timesteps. This smoothness persists under numerical discretization, providing a foundation for feature caching strategies.*

**Diffusion Transformer Architecture.** The Diffusion Transformer (DiT) follows a hierarchical structure, $\mathcal{G} = g_1 \circ g_2 \circ \cdots \circ g_L$, where each block $g_l = \mathcal{F}_{SA}^l \circ \mathcal{F}_{CA}^l \circ \mathcal{F}_{MLP}^l$ consists of self-attention (SA), cross-attention (CA), and multilayer perceptron (MLP) components. The superscript $l \in \{1, 2, ..., L\}$ denotes the layer index. In DiT, both attention mechanisms and MLP components evolve over time. At each timestep, $\mathcal{F}_{SA}^l$, $\mathcal{F}_{CA}^l$, and $\mathcal{F}_{MLP}^l$ dynamically adjust to accommodate varying noise levels during image generation. The input, $\mathbf{x}_t = \{x_i\}_{i=1}^{H \times W}$, is represented as a sequence of tokens corresponding to image patches. Each block integrates residual connections, expressed as $\mathcal{F}(\mathbf{x}) = \mathbf{x} + \text{AdaLN} \circ f(\mathbf{x})$, where $f(\mathbf{x})$ corresponds to MLP or attention layers, and AdaLN denotes adaptive layer normalization, ensuring consistent information flow.

**Naïve Feature Caching for Diffusion Models.** Recent acceleration methods employ *Naïve Feature Caching Strategies* in diffusion models by directly reusing computed features across adjacent timesteps. Specifically, given timesteps $\{t, t-1, \ldots, t-(N-1)\}$, features computed at timestep $t$ are cached as $\mathcal{C}(x_t^l) := \{\mathcal{F}(x_t^l)\}$. These cached features are then directly reused for subsequent steps: $\mathcal{F}(x_{t-k}^l) := \mathcal{F}(x_t^l)$, where $k \in 1, \ldots, N-1$. While this approach achieves a theoretical $(N-1)$-fold speedup by eliminating redundant computations, it suffers from exponential error accumulation as $N$ increases due to neglecting the temporal dynamics of features.

### 3.2. TaylorSeer

*TaylorSeer* introduces Taylor series-based predictive caching to mitigate error accumulation in conventional caching approaches. By leveraging the continuous nature of feature trajectory, we develop a predictive caching method that accurately estimates intermediate features across timesteps.

**Linear Prediction Method.** To overcome the limitations of *naïve caching*, we propose a *linear prediction strategy* that extends beyond direct feature reuse. Instead of directly copying features, we cache both feature values and their temporal differences: $\mathcal{C}(x_t^l) := \{\mathcal{F}(x_t^l), \Delta\mathcal{F}(x_t^l)\}$. This allows us to predict feature trajectories at timestep $t-k$ using the following formulation:

$$\mathcal{F}_{\text{pred}}(x_{t-k}^l) = \mathcal{F}(x_t^l) + \frac{\mathcal{F}(x_t^l) - \mathcal{F}(x_{t+N}^l)}{N}k \quad (4)$$

where $\mathcal{F}(x_{t+N}^l)$ and $\mathcal{F}(x_t^l)$ denote features at fully activated timesteps. This first-order approximation captures the linear trend of feature trajectories, significantly improving accuracy over direct feature reuse.

**Higher-Order Prediction via Taylor Expansion.** To further improve the accuracy of feature prediction, we extend our forecasting method by incorporating Taylor's theorem. This approach leverages higher-order finite difference approximations to capture the temporal dynamics of features, reducing cumulative prediction errors while maintaining computational efficiency. At timestep $t$, we define a cache storing the feature and its $m$-th order finite differences:

$$\mathcal{C}(x_t^l) := \{\mathcal{F}(x_t^l), \Delta\mathcal{F}(x_t^l), ..., \Delta^m\mathcal{F}(x_t^l)\} \quad (5)$$

where $\Delta^i\mathcal{F}(x_t^l)$ represents the $i$-th order finite difference, approximating the temporal dynamics of $\mathcal{F}(x_t^l)$ around timestep $t$. For an $(m+1)$-times differentiable feature function $\mathcal{F}(x_t^l)$, the feature at timestep $t-k$ can be expressed using a Taylor series expansion:

$$\mathcal{F}(x_{t-k}^l) = \mathcal{F}(x_t^l) + \sum_{i=1}^{m} \frac{\mathcal{F}^{(i)}(x_t^l)}{i!}(-k)^i + R_{m+1} \quad (6)$$

where $R_{m+1}$ represents the remainder term. To avoid explicit computation of higher-order derivatives, we approximate them using finite differences. The $i$-th order forward finite difference is defined recursively as:

$$\Delta^i\mathcal{F}(x_t^l) = \Delta(\Delta^{i-1}\mathcal{F}(x_t^l)) = \Delta^{i-1}\mathcal{F}(x_{t+N}^l) - \Delta^{i-1}\mathcal{F}(x_t^l) \quad (7)$$

with the base case $\Delta^0\mathcal{F}(x_t^l) = \mathcal{F}(x_t^l)$. Equation (7) can be further expanded into a binomial form:

$$\Delta^i\mathcal{F}(x_t^l) = \sum_{j=0}^{i} (-1)^{i-j} \binom{i}{j} \mathcal{F}(x_{t+jN}^l) \quad (8)$$

4

It can be shown that the $i$-th order finite difference approximates the $i$-th derivative scaled by $N^i$:

$$\Delta^i \mathcal{F}(x_t^l) \approx N^i \mathcal{F}^{(i)}(x_t^l) \quad (9)$$

Substituting this approximation into the Taylor expansion and accounting for the scaling factor, we derive the $m$-th order prediction formula:

$$\mathcal{F}_{\text{pred},m}(x_{t-k}^l) = \mathcal{F}(x_t^l) + \sum_{i=1}^{m} \frac{\Delta^i \mathcal{F}(x_t^l)}{i! \cdot N^i}(-k)^i \quad (10)$$

This formulation requires only $(m + 1)$ fully computed timesteps $\{t + mN, \ldots, t + N, t\}$ to predict features at intermediate timesteps, achieving an optimal balance between efficiency and accuracy.

Our method establishes a principled transition from directly feature caching to predictive forecasting by leveraging the mathematical foundation of Taylor series approximation. Instead of merely reusing features, we model their temporal trajectory over time, enabling accurate estimation of intermediate representations and capturing both short-term and long-term dynamics. This systematic formulation unifies various prediction strategies—ranging from simple caching to higher-order forecasting—within a single coherent methodology, providing flexibility and robustness across different temporal scales and scenarios.

- **Directly Caching** ($m = 0$): Reduces to *Naïve Feature Caching*, reusing features without temporal modeling. Simple but less accurate for dynamic feature changes.
- **Short-Term Forecasting** ($m = 1$): Uses first-order finite differences for linear prediction, suitable for short-term changes but limited for complex dynamics.
- **Long-Range Forecasting** ($m \geq 2$): Leverages higher-order finite differences to model nonlinear trajectories, reducing errors and capturing long-term temporal patterns.

Unlike Naïve Feature Caching, *TaylorSeer* transforms simply feature reuse into a predictive process, explicitly modeling feature trajectories through high-order finite differences. This key transition from directly caching to forecasting not only significantly improves prediction accuracy but also enables efficient long-range inference.. As a result, *TaylorSeer* provides a robust and scalable solution for accelerating diffusion model inference, maintaining high generation quality even at large acceleration ratios.

**Error Bounds Analysis.** For a feature function $\mathcal{F}(x_t^l)$ that is $(m+1)$-times differentiable on $[t-k, t]$, we define the prediction error as $E_m(k) = \|\mathcal{F}_{\text{pred},m}(x_{t-k}^l) - \mathcal{F}(x_{t-k}^l)\|$. By Taylor's remainder theorem, this error is bounded by:

$$E_m(k) \leq \frac{M_{m+1}}{(m+1)!}|k|^{m+1} \quad (11)$$

where $M_{m+1} = \sup_{\xi \in [t-k,t]} \|\mathcal{F}^{(m+1)}(x_\xi^l)\|$. Since our method employs finite difference approximations rather

than exact derivatives, the complete error bound incorporates additional terms:

$$E_m(k) \leq \frac{M_{m+1}}{(m+1)!}|k|^{m+1} + \sum_{i=1}^{m} \frac{C_i}{i!}|k|^i|N|^{i-1} \quad (12)$$

where constants $C_i$ relate to finite difference approximation errors. This analysis reveals a fundamental trade-off: higher-order predictions ($m$) effectively reduce the primary error term but introduce additional errors scaling with the sampling interval $N$. For diffusion models with smooth feature trajectories (small $M_{m+1}$), our method achieves optimal accuracy by balancing prediction order and caching interval, particularly effective when $|k| < |N|$ where the Taylor approximation dominates the error bound.

## 4. Experiments

### 4.1. Experiment Settings

**Model Configurations.** The experiments are conducted on three state-of-the-art visual generative models: the text-to-image generation model FLUX.1-dev [16], text-to-video generation model HunyuanVideo [43], and the class-conditional image generation model DiT-XL/2 [28].
**FLUX.1-dev**[16] predominantly employs the Rectified Flow [22] sampling method with 50 steps as the standard configuration. All experimental evaluations of FLUX.1-dev were conducted on NVIDIA H800 GPUs.
**HunyuanVideo** [20, 43] was evaluated on the Hunyuan-Large architecture, utilizing the standard 50-step inference protocol as the baseline while preserving all default sampling parameters for rigorous experimental consistency. Extensive performance benchmarks were systematically conducted using NVIDIA H20 96GB GPUs for detailed latency assessment and NVIDIA H100 80GB GPUs for comprehensive inference operations.
**DiT-XL/2** [28] adopts a 50-step DDIM [40] sampling strategy to ensure consistency with other models. All models incorporate a unified forced activation period $\mathcal{N}$, while $\mathcal{O}$ represents the order of the Taylor expansion, optimizing computational efficiency and overall model performance. Experiments on DiT-XL/2 are conducted on NVIDIA A800 80GB GPUs. *For more detailed model configurations, please refer to the Supplementary Material.*

**Evaluation and Metrics.** For the text-to-image generation task, we perform inference on 200 DrawBench [35] prompts to generate images with a resolution of $1000 \times 1000$. We then evaluate the generated samples using Image Reward [44] and CLIP Score [9] as key metrics to assess image quality and text alignment. For the text-to-video generation task, we leverage the VBench [13] evaluation framework, utilizing its 946 benchmark prompts. For each prompt, we generate five samples with different random seeds, totaling 4,730 videos. We then systematically

Table 1. **Quantitative comparison in text-to-image generation** for FLUX on Image Reward.

| Method FLUX.1[16] | Efficient Attention [7] | Acceleration | | | | Image Reward ↑ DrawBench | CLIP↑ Score |
|---|---|---|---|---|---|---|---|
| | | Latency(s) ↓ | Speed ↑ | FLOPs(T) ↓ | Speed ↑ | | |
| **[dev]: 50 steps** | ✔ | 17.20 | 1.00× | 3719.50 | 1.00× | 0.9898 | 19.604 |
| 60% **steps** | ✔ | 10.49 | 1.64× | 2231.70 | 1.67× | 0.9739 | 19.526 |
| $\Delta$-DiT ($\mathcal{N}=2$) † | ✔ | 11.87 | 1.45× | 2480.01 | 1.50× | 0.9316 | 19.350 |
| 50% **steps** † | ✔ | 8.80 | 1.95× | 1859.75 | 2.00× | 0.9429 | 19.325 |
| 40% **steps** † | ✔ | 7.11 | 2.42× | 1487.80 | 2.62× | 0.9317 | 19.027 |
| 34% **steps** † | ✔ | 6.09 | 2.82× | 1264.63 | 3.13× | 0.9346 | 18.904 |
| $\Delta$-DiT ($\mathcal{N}=3$) † | ✔ | 8.81 | 1.95× | 1686.76 | 2.21× | 0.8561 | 18.833 |
| **FORA** ($\mathcal{N}=3$) † [37] | ✔ | 7.08 | 2.43× | 1320.07 | 2.82× | 0.9227 | 18.950 |
| **ToCa** ($\mathcal{N}=5$)[53] | ✘ | 10.80 | 1.59× | 1126.76 | 3.30× | 0.9731 | 19.030 |
| **DuCa**($\mathcal{N}=5$) [54] | ✔ | 5.88 | 2.93× | 1078.34 | 3.45× | 0.9896 | **19.595** |
| **TaylorSeer** ($\mathcal{N}=5, O=2$) | ✔ | 5.82 | 2.96× | **893.54** | **4.16×** | **1.0296** | 19.437 |
| **FORA** ($\mathcal{N}=4$) † [37] | ✔ | 5.43 | 3.17× | 967.91 | 3.84× | 0.8675 | 18.560 |
| **ToCa** ($\mathcal{N}=8$) † [53] | ✘ | 8.47 | 2.03× | 784.54 | 4.74× | 0.9086 | 18.380 |
| **DuCa** ($\mathcal{N}=6$) † [54] | ✔ | 4.89 | 3.52× | 816.55 | 4.56× | 0.9470 | 19.082 |
| **TaylorSeer** ($\mathcal{N}=6, O=1$) | ✔ | 4.87 | 3.53× | **744.81** | **4.99×** | 0.9953 | **19.637** |
| **TaylorSeer** ($\mathcal{N}=6, O=2$) | ✔ | 5.19 | 3.31× | 744.81 | 4.99× | **1.0039** | 19.427 |

- † Methods exhibit significant degradation in Image Reward, leading to severe deterioration in image quality.

Table 2. **Quantitative comparison in text-to-video generation** for HunyuanVideo on VBench.

| Method HunyuanVideo[43] | Efficient Attention [7] | Acceleration | | | | VBench ↑ Score(%) |
|---|---|---|---|---|---|---|
| | | Latency(s) ↓ | Speed ↑ | FLOPs(T) ↓ | Speed ↑ | |
| **Original: 50 steps** | ✔ | 318.24 | 1.00× | 29773.0 | 1.00× | 80.66 |
| 22% **steps** | ✔ | 70.34 | 4.52× | 6550.1 | 4.55× | 78.74 |
| **FORA** [37] | ✔ | 67.19 | 4.74× | 5960.4 | 5.00× | 78.83 |
| **ToCa** [53] | ✘ | 77.82* | 4.09× | 7006.2 | 4.25× | 78.86 |
| **DuCa** [54] | ✔ | 71.10 | 4.48× | 6483.2 | 4.62× | 78.72 |
| **TaylorSeer** ($\mathcal{N}=5, O=1$) | ✔ | 68.42 | 4.65× | 5960.4 | 5.00× | **79.93** |
| **TaylorSeer** ($\mathcal{N}=6, O=1$) | ✔ | **61.99** | **5.13×** | **5359.1** | **5.56×** | 79.78 |

assess the generated results based on 16 core evaluation dimensions defined by the VBench framework to provide a comprehensive evaluation of the model's performance. For the class-conditional image generation task, we uniformly sample from 1,000 ImageNet [33] categories, generating 50,000 images with a resolution of $256 \times 256$. We use FID-50k [10] as the primary evaluation metric, complemented by sFID and Inception Score for robust evaluation.

### 4.2. Text-to-Image Generation

**Quantitative Study.** We compared *TaylorSeer* with existing methods. While DuCa [54] ($\mathcal{N}=5$) achieves a 3.45× FLOPs speedup with an Image Reward of 0.9896, and ToCa [53] ($\mathcal{N}=5$) offers 3.30× acceleration with reduced quality (0.9731), *TaylorSeer* ($\mathcal{N}=6$, $\mathcal{O}=1$) significantly outperforms them. At **4.99× acceleration**, it maintains superior Image Reward (0.9953) and CLIP Score (19.637). Notably, when pushing acceleration further, baseline methods suffer severe quality degradation: DuCa ($\mathcal{N}=6$) drops to 0.9470 Image Reward at 4.56× acceleration, while ToCa ($\mathcal{N}=8$) plummets to 0.9086 at 4.74× acceleration. In

contrast, *TaylorSeer* ($\mathcal{N}=5$, $\mathcal{O}=2$) sustains 1.0296 Image Reward even at **4.16× acceleration**, demonstrating unmatched efficiency-fidelity balance.

**Qualitative Study.** The qualitative results highlight *TaylorSeer*'s superior ability to preserve image quality while accelerating computation. In text generation tasks, such as *"The word 'START' written on a street surface,"* *TaylorSeer* accurately retains the textual elements, whereas methods like FORA and ToCa lose critical details. In Renaissance-style portrait generation, *TaylorSeer* maintains consistent image fidelity, while other methods exhibit notable *color inaccuracies* and *missing objects* across test cases. This shows *TaylorSeer* balances efficiency and quality, particularly in tasks needing fine detail preservation.

### 4.3. Text-to-Video Generation

**Quantitative Study.** On HunyuanVideo, *TaylorSeer* with $\mathcal{N}=5$ and $\mathcal{O}=1$ reduces inference latency to 68 seconds and computational cost to 5960.4TFLOPs (5.00× speedup),

Table 3. **Quantitative comparison on class-to-image generation** on ImageNet with DiT-XL/2.

| Method | Efficient Attention | Latency(s) ↓ | FLOPs(T) ↓ | Speed ↑ | FID ↓ | sFID ↓ | Inception Score ↑ |
|---|---|---|---|---|---|---|---|
| **DDIM-50 steps** | ✔ | 0.428 | 23.74 | 1.00× | 2.32 | 4.32 | 241.25 |
| **DDIM-25 steps** | ✔ | 0.230 | 11.87 | 2.00× | 3.18 | 4.74 | 232.01 |
| **Δ-DiT**($\mathcal{N}=2$) | ✔ | 0.246 | 18.04 | 1.31× | 2.69 | 4.67 | 225.99 |
| **Δ-DiT**($\mathcal{N}=3$) † | ✔ | 0.173 | 16.14 | 1.47× | 3.75 | 5.70 | 207.57 |
| **DDIM-20 steps** | ✔ | 0.191 | 9.49 | 2.50× | 3.81 | 5.15 | 221.43 |
| **FORA** ($\mathcal{N}=3$) | ✔ | 0.197 | 8.58 | 2.77× | 3.55 | 6.36 | 229.02 |
| **ToCa** ($\mathcal{N}=3$) | ✘ | 0.216 | 10.23 | 2.32× | 2.87 | 4.76 | 235.21 |
| **DuCa** ($\mathcal{N}=3$) | ✔ | 0.208 | 9.58 | 2.48× | 2.88 | **4.66** | 233.37 |
| **TaylorSeer** ($\mathcal{N}=3, O=3$) | ✔ | 0.248 | **8.56** | **2.77×** | 2.34 | 4.69 | **238.42** |
| **DDIM-15 steps** † | ✔ | 0.152 | 7.12 | 3.33× | 5.17 | 6.11 | 206.33 |
| **FORA** ($\mathcal{N}=4$) † | ✔ | 0.169 | 6.66 | 3.56× | 4.75 | 8.43 | 213.72 |
| **ToCa** ($\mathcal{N}=5$) † | ✘ | 0.173 | 6.77 | 3.51× | 6.20 | 7.17 | 200.04 |
| **DuCa** ($\mathcal{N}=4$)† | ✔ | 0.144 | 7.61 | 3.11× | 3.42 | **4.94** | 225.19 |
| **TaylorSeer** ($\mathcal{N}=4, O=4$) | ✔ | 0.220 | **6.66** | **3.56×** | 2.49 | 5.19 | **235.83** |
| **DDIM-12 steps** † | ✔ | 0.128 | 5.70 | 4.17× | 7.80 | 8.03 | 184.50 |
| **DDIM-10 steps** † | ✔ | 0.115 | 4.75 | 5.00× | 12.15 | 11.33 | 159.13 |
| **FORA** ($\mathcal{N}=5$) † | ✔ | 0.149 | 5.24 | 4.53× | 6.58 | 11.29 | 193.01 |
| **ToCa** ($\mathcal{N}=6$) † | ✘ | 0.163 | 6.34 | 3.75× | 6.55 | 7.10 | 189.53 |
| **DuCa** ($\mathcal{N}=5$)† | ✔ | 0.154 | 6.27 | 3.78× | 6.06 | 6.72 | 198.46 |
| **TaylorSeer** ($\mathcal{N}=5, O=3$) | ✔ | 0.180 | **5.24** | **4.53×** | 2.65 | 5.36 | **231.59** |

- † Methods exhibit significant degradation in FID, leading to severe deterioration in image quality.
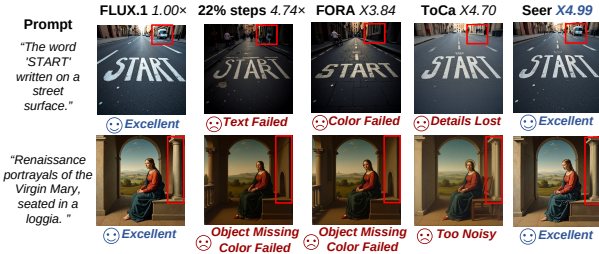


Figure 4. Detailed visualization results for different acceleration methods on FLUX.1-dev. Other methods exhibit issues such as text failure, color distortion, and missing details, whereas TaylorSeer achieves the best quality and acceleration.

achieving a 79.93% VBench score that outperforms both ToCa and DuCa. With $\mathcal{N}$=6 and $\mathcal{O}$=1, performance improves further to $5.56\times$ speedup while maintaining a 79.78% score. These results show that *TaylorSeer* preserves temporal consistency and reduces computational demands, making high-quality video synthesis more efficient.

**Qualitative Study.** The qualitative results highlight *TaylorSeer*'s ability to preserve video quality while accelerating computation. In the *"a fire hydrant and a stop sign"* scenario, *TaylorSeer* accurately generates the correct text, while methods like HunyuanVideo introduce a spelling error in *"STOP"*, affecting semantic accuracy. In the *"a donut and a suitcase"* case, *TaylorSeer* maintains exceptional visual consistency, whereas other methods fail to generate the suitcase entirely. In the *"a motorcycle turning a corner"* case, *TaylorSeer* preserves motion blur details

and trajectory smoothness, leading to more realistic video. Across these diverse scenarios, *TaylorSeer* consistently outperforms other methods in fidelity and quality

## 4.4. Class-Conditional Image Generation

We compared *TaylorSeer* with methods such as ToCa [53], FORA [37], DuCa [54], and reduced DDIM steps on DiT-XL/2 [28], demonstrating that *TaylorSeer* significantly outperforms the others in both acceleration ratio and generation quality. *TaylorSeer*($\mathcal{N}$=3, $\mathcal{O}$=3) achieves an FID-50k of 2.34 while providing a **2.77× acceleration**, showing that higher-order Taylor expansions preserve feature quality at high acceleration ratios without incurring additional computational overhead. At a **4.53× acceleration**, our method ($\mathcal{N}$=5, $\mathcal{O}$=3) maintains an FID of 2.65, outperforming state-of-the-art models such as ToCa [53] and DuCa [54]. Notably, as the acceleration ratio increases beyond $3.5\times$, other methods like FORA, ToCa, and DuCa exhibit significant degradation in FID, leading to severe deterioration in image quality, while *TaylorSeer* consistently maintains superior performance. This robustness under extreme acceleration conditions further validates the effectiveness of our approach in predicting future features, which minimizes error accumulation compared to cache-then-reuse strategies that struggle to adapt to higher acceleration demands.

## 4.5. Ablation Studies

We conduct ablation experiments on DiT-XL/2 [28] to evaluate *TaylorSeer*, focusing on the impact of the interval pa-
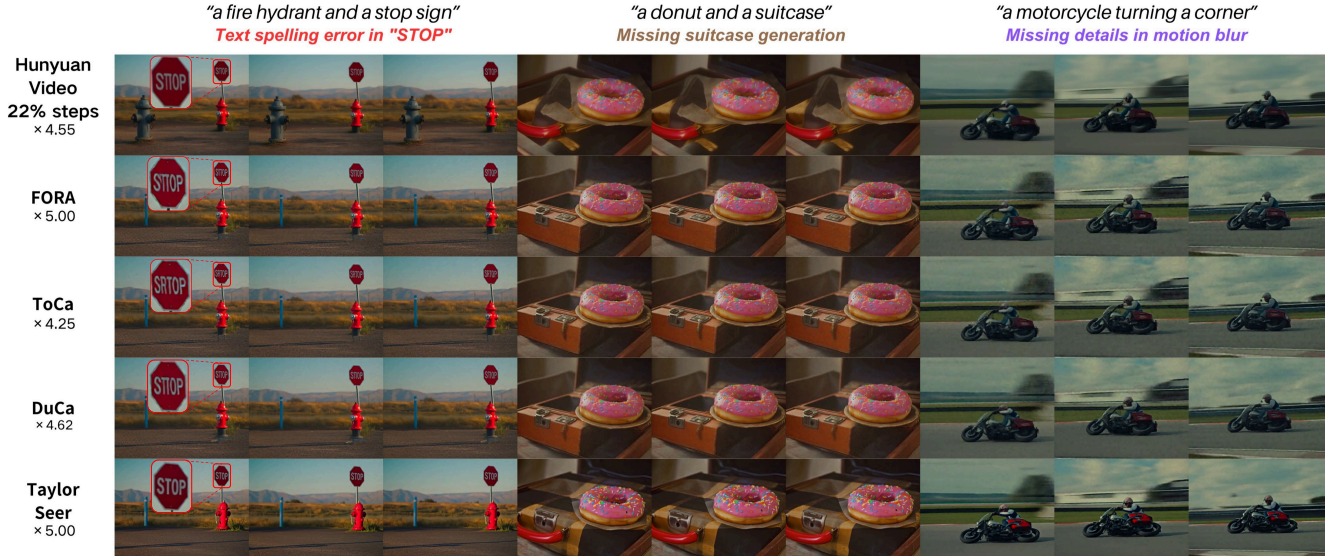
Figure 5. Visualization of different acceleration methods on HunyuanVideo. While achieving higher acceleration ratios, other methods exhibit issues such as *text errors*, *missing content*, and *motion detail loss*. In contrast, our method demonstrates superior performance, maintaining high-quality generation without these problems.
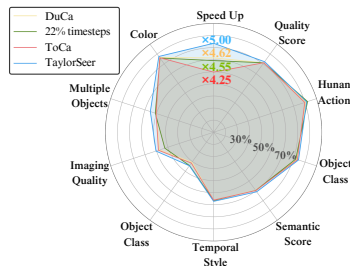


Figure 6. VBench metrics and acceleration ratios of *TaylorSeer* and other methods.
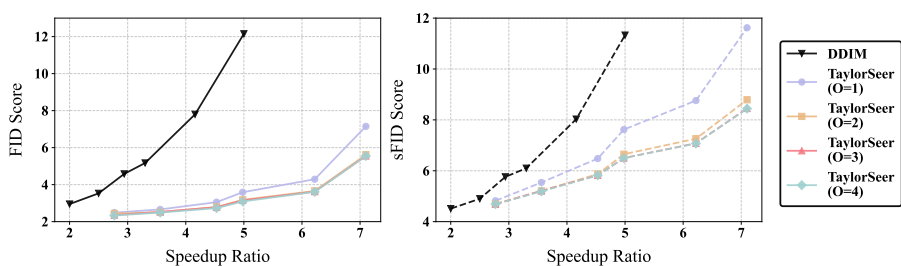
Figure 7. Comparison of 1st–4th order Taylor expansions for caching-based feature prediction. Higher-order expansions exhibit similar convergence and consistently outperform the 1st-order expansion at lower FLOPs, demonstrating the advantages of high-order approximation.

rameter $\mathcal{N}$ and the Taylor expansion order $\mathcal{O}$ on computational efficiency and generation quality. Results indicate that when $\mathcal{O}$=0 (direct feature reuse), performance degrades rapidly as $\mathcal{N}$ increases. Introducing higher-order Taylor expansions significantly improves generation quality, especially in long-interval scenarios. Even a first-order approximation ($\mathcal{O}$=1) substantially enhances performance, reducing FID from non-generation to 4.29 at $\mathcal{N}$=7. Higher-order expansions ($\mathcal{O} \geq 2$) further improve generation quality, capturing nonlinear feature variations effectively. In high-acceleration settings ($\mathcal{N}$=5 or 6), third-order expansion maintains low FID (2.78–3.15), outperforming direct reuse. While $\mathcal{O} \geq 3$ continues to refine results, improvements saturate beyond the third order, as observed in both sFID and FID trends.

Overall, *TaylorSeer* demonstrates that Taylor expansion-based forecasting effectively balances efficiency and quality, particularly in high-acceleration scenarios. Higher-order expansions ($\mathcal{O} \geq 3$) enable high-quality generation with reduced computational cost, making the method suit-

able for real-time or resource-constrained applications. *Detailed results are provided in the supplementary materials.*

## 5. Conclusion

Traditional feature caching methods follow the paradigm of *"cache-then-reuse"*, which directly reuses the features stored at previous timesteps in the following timesteps, and suffer from serve drop in generation quality at high acceleration ratios. In this paper, motivated by the surprisingly stable trajectory of features at different timesteps, we propose *"cache-then-forecast"*, which formulates feature caching as a sequential prediction problem and solves it with Taylor expansion. The high-order Taylor expansions in our method can capture complex feature trajectories more rigorously, making it maintain the generation quality in high-acceleration ratios, where all previous caching methods totally fail. Experiments across architectures including DiT, FLUX, and HunyuanVideo demonstrate significant acceleration ($4.53\times$–$5.56\times$) without quality degradation. We hope *TaylorSeer* can move the paradigm of feature caching methods from reusing to forecasting.

# References

[1] 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. 1

[2] Daniel Bolya and Judy Hoffman. Token merging for fast stable diffusion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4599–4603, 2023. 3

[3] 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, 2024. 2

[4] 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. In *International Conference on Learning Representations*, 2024. 2

[5] Pengtao Chen, Mingzhu Shen, Peng Ye, Jianjian Cao, Chongjun Tu, Christos-Savvas Bouganis, Yiren Zhao, and Tao Chen. $\delta$-dit: A training-free acceleration method tailored for diffusion transformers. *arXiv preprint arXiv:2406.01125*, 2024. 1, 3

[6] Xinle Cheng, Zhuoming Chen, and Zhihao Jia. Cat pruning: Cluster-aware token pruning for text-to-image diffusion models, 2025. 3

[7] Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher R'e. Flashattention: Fast and memory-efficient exact attention with io-awareness. *ArXiv*, abs/2205.14135, 2022. 6

[8] Gongfan Fang, Xinyin Ma, and Xinchao Wang. Structural pruning for diffusion models. *arXiv preprint arXiv:2305.10924*, 2023. 3

[9] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. CLIPScore: A Reference-free Evaluation Metric for Image Captioning, 2022. arXiv:2104.08718 [cs]. 5

[10] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium, 2018. arXiv:1706.08500 [cs]. 6

[11] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models, 2020. arXiv:2006.11239 [cs]. 1

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

[13] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, Yaohui Wang, Xinyuan Chen, Limin Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. VBench: Comprehensive Benchmark Suite for Video Generative Models, 2023. arXiv:2311.17982 [cs]. 5

[14] Minchul Kim, Shangqian Gao, Yen-Chang Hsu, Yilin Shen, and Hongxia Jin. Token fusion: Bridging the gap between token pruning and token merging. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1383–1392, 2024. 3

[15] Sungbin Kim, Hyunwuk Lee, Wonho Cho, Mincheol Park, and Won Woo Ro. Ditto: Accelerating diffusion model via temporal value similarity. In *Proceedings of the 2025 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 2025. 3

[16] Black Forest Labs. Flux. https://github.com/black-forest-labs/flux, 2024. 5, 6, 1

[17] Senmao Li, Taihang Hu, Fahad Shahbaz Khan, Linxuan Li, Shiqi Yang, Yaxing Wang, Ming-Ming Cheng, and Jian Yang. Faster diffusion: Rethinking the role of unet encoder in diffusion models. *arXiv preprint arXiv:2312.09608*, 2023. 3

[18] Xiuyu Li, Yijiang Liu, Long Lian, Huanrui Yang, Zhen Dong, Daniel Kang, Shanghang Zhang, and Kurt Keutzer. Q-diffusion: Quantizing diffusion models. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 17489–17499, 2023. 3

[19] Yanyu Li, Huan Wang, Qing Jin, Ju Hu, Pavlo Chemerys, Yun Fu, Yanzhi Wang, Sergey Tulyakov, and Jian Ren. Snapfusion: Text-to-image diffusion model on mobile devices within two seconds. *Advances in Neural Information Processing Systems*, 36, 2024. 3

[20] Zhimin Li, Jianwei Zhang, and and others Lin. Hunyuan-DiT: A powerful multi-resolution diffusion transformer with fine-grained chinese understanding. 5, 1

[21] Feng Liu, Shiwei Zhang, Xiaofeng Wang, Yujie Wei, Haonan Qiu, Yuzhong Zhao, Yingya Zhang, Qixiang Ye, and Fang Wan. Timestep embedding tells: It's time to cache for video diffusion model, 2024. 3

[22] Xingchao Liu, Chengyue Gong, et al. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *The Eleventh International Conference on Learning Representations*, 2023. 2, 5

[23] Ziming Liu, Yifan Yang, Chengruidong Zhang, Yiqi Zhang, Lili Qiu, Yang You, and Yuqing Yang. Region-adaptive sampling for diffusion transformers, 2025. 3

[24] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787, 2022. 2

[25] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022. 2

[26] Xinyin Ma, Gongfan Fang, and Xinchao Wang. Deepcache: Accelerating diffusion models for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15762–15772, 2024. 1, 3

[27] Chenlin Meng, Ruiqi Gao, Diederik P Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. In *NeurIPS 2022 Workshop on Score-Based Methods*, 2022. 2

[28] William Peebles and Saining Xie. Scalable Diffusion Models with Transformers, 2023. arXiv:2212.09748 [cs]. 1, 5, 7

[29] 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. 2

[30] Junxiang Qiu, Shuo Wang, Jinda Lu, Lin Liu, Houcheng Jiang, and Yanbin Hao. Accelerating diffusion transformer via error-optimized cache, 2025. 3

[31] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-Resolution Image Synthesis with Latent Diffusion Models, 2022. arXiv:2112.10752 [cs]. 1

[32] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015. 2

[33] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge, 2015. arXiv:1409.0575 [cs]. 6

[34] Omid Saghatchian, Atiyeh Gh. Moghadam, and Ahmad Nickabadi. Cached adaptive token merging: Dynamic token reduction and redundant computation elimination in diffusion model, 2025. 3

[35] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding. 5

[36] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022. 2

[37] Pratheba Selvaraju, Tianyu Ding, Tianyi Chen, Ilya Zharkov, and Luming Liang. Fora: Fast-forward caching in diffusion transformer acceleration. *arXiv preprint arXiv:2407.01425*, 2024. 1, 3, 6, 7

[38] Yuzhang Shang, Zhihang Yuan, Bin Xie, Bingzhe Wu, and Yan Yan. Post-training quantization on diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1972–1981, 2023. 3

[39] 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. 2

[40] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021. 2, 5

[41] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. In *International Conference on Machine Learning*, pages 32211–32252. PMLR, 2023. 2

[42] Wenzhang Sun, Qirui Hou, Donglin Di, Jiahui Yang, Yongjia Ma, and Jianxun Cui. Unicp: A unified caching and pruning framework for efficient video generation, 2025. 3

[43] Xingwu Sun, Yanfeng Chen, Huang, et al. Hunyuan-large: An open-source MoE model with 52 billion activated parameters by tencent. 5, 6, 1

[44] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. ImageReward: Learning and Evaluating Human Preferences for Text-to-Image Generation, 2023. arXiv:2304.05977 [cs]. 5

[45] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, Da Yin, Xiaotao Gu, Yuxuan.Zhang, Weihan Wang, Yean Cheng, Bin Xu, Yuxiao Dong, and Jie Tang. Cogvideox: Text-to-video diffusion models with an expert transformer. In *The Thirteenth International Conference on Learning Representations*, 2025. 2

[46] Zhihang Yuan, Hanling Zhang, Lu Pu, Xuefei Ning, Linfeng Zhang, Tianchen Zhao, Shengen Yan, Guohao Dai, and Yu Wang. DiTFastattn: Attention compression for diffusion transformer models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 1, 3

[47] Evelyn Zhang, Bang Xiao, Jiayi Tang, Qianli Ma, Chang Zou, Xuefei Ning, Xuming Hu, and Linfeng Zhang. Token pruning for caching better: 9 times acceleration on stable diffusion for free, 2024. 3

[48] Evelyn Zhang, Jiayi Tang, Xuefei Ning, and Linfeng Zhang. Training-free and hardware-friendly acceleration for diffusion models via similarity-based token pruning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2025. 3

[49] Xuanlei Zhao, Xiaolong Jin, Kai Wang, and Yang You. Real-time video generation with pyramid attention broadcast. *arXiv preprint arXiv:2408.12588*, 2024. 1

[50] Kaiwen Zheng, Cheng Lu, Jianfei Chen, and Jun Zhu. DPM-solver-v3: Improved diffusion ODE solver with empirical model statistics. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 2

[51] 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, 2024. 2

[52] Haowei Zhu, Dehua Tang, Ji Liu, Mingjie Lu, Jintu Zheng, Jinzhang Peng, Dong Li, Yu Wang, Fan Jiang, Lu Tian, Spandan Tiwari, Ashish Sirasao, Jun-Hai Yong, Bin Wang, and Emad Barsoum. Dip-go: A diffusion pruner via few-step gradient optimization, 2024. 3

[53] Chang Zou, Xuyang Liu, Ting Liu, Siteng Huang, and Linfeng Zhang. Accelerating diffusion transformers with token-wise feature caching. *arXiv preprint arXiv:2410.05317*, 2024. 1, 3, 6, 7

[54] Chang Zou, Evelyn Zhang, Runlin Guo, Haohang Xu, Conghui He, Xuming Hu, and Linfeng Zhang. Accelerating diffusion transformers with dual feature caching, 2024. 1, 3, 6, 7

# From Reusing to Forecasting: Accelerating Diffusion Models with TaylorSeers

## Supplementary Material

## 6. Experimental Details

In this section, more details of the experiments are provided.

### 6.1. Model Configuration

As mentioned in 4.1.1 , experiments on three models from different tasks, FLUX [16] for text-to-image generation, HunyuanVideo [20, 43] for text-to-video generation, and DiT [28] for class-conditional image generation, are presented. In this section, a more detailed hyperparameter configuration scheme is provided.

- **FLUX**: The FORA [37] method employs a uniform activation interval with $\mathcal{N}$=3. The ToCa [53] method uses $\mathcal{N}$=4 with a caching ratio of 90%, adopting a non-uniform activation interval, with sparse activations at the beginning and dense activations towards the end, utilizing an attention-based token selection method. The DuCa [54] method sets conservative caching steps on even-numbered steps following fresh steps, while aggressive caching steps are set on odd-numbered steps. The activation interval and caching ratio are consistent with the ToCa method, also using a non-uniform activation interval and employing the attention-based token selection method.
- **HunyuanVideo**: The FORA [37] method utilizes an activation interval of $\mathcal{N}$=5, whereas both the ToCa [53] and DuCa methods employ $\mathcal{N}$=4 with a caching ratio of 90%. For each activation step (complete computational step), aggressive caching is applied to odd-numbered steps, and conservative caching is applied to even-numbered steps. Due to memory limitations that result in an "out of memory" error when using a non-uniform activation scheme, the ToCa method in HunyuanVideo is configured with a uniform activation interval, as indicated by a $^*$ in the corresponding table.
- **DiT**: The FORA [37] method uses a uniform activation interval with $\mathcal{N}$=3. The ToCa [53] method also uses $\mathcal{N}$=3 with an average caching ratio of $R = 95\%$, employing a non-uniform activation interval, with sparse activations at the beginning and dense activations towards the end, utilizing the attention-based token selection method. The DuCa [54] method sets conservative caching steps on even-numbered steps following fresh steps, while aggressive caching steps are set on odd-numbered steps. The activation interval and caching ratio are consistent with the ToCa method, also using a non-uniform activation interval with sparse-to-dense activation and employing the attention-based token selection method. $\Delta$-DiT adopts a layer-skipping strategy

where, in the early stages (49-25 steps), layers 14-27 are skipped, and in the later stages (24-0 steps), layers 0-13 are skipped.

## 7. Supplementary Results for Ablation Studies

We conduct ablation experival parameter $\mathcal{N}$ and the Taylor expansion order $\mathcal{O}$ on computatments on DiT-XL/2 [28] to evaluate *TaylorSeer*, focusing on the impact of the interional efficiency and generation quality. The results demonstrate the importance of these design choices in balancing performance and speed.

Table 4. **Ablation Study with Different Configurations** on ImageNet with DiT-XL/2.

| Configuration | FLOPs(T)↓ | Speed↑ | sFID↓ | FID↓ |
|---|---|---|---|---|
| ($\mathcal{N}$=3, $\mathcal{O}$=0) | 8.56 | 2.77× | 6.36 | 3.55 |
| ($\mathcal{N}$=4, $\mathcal{O}$=0) | 6.66 | 3.56× | 8.43 | 4.75 |
| ($\mathcal{N}$=5, $\mathcal{O}$=0) | 5.24 | 4.53× | 11.29 | 6.58 |
| ($\mathcal{N}$=6, $\mathcal{O}$=0) | 4.76 | 4.98× | 14.84 | 9.24 |
| ($\mathcal{N}$=7, $\mathcal{O}$=0) | 3.82 | 6.22× | 18.57 | 12.67 |
| ($\mathcal{N}$=3, $\mathcal{O}$=1) | 8.56 | 2.77× | 4.82 | 2.49 |
| ($\mathcal{N}$=4, $\mathcal{O}$=1) | 6.66 | 3.56× | 5.54 | 2.66 |
| ($\mathcal{N}$=5, $\mathcal{O}$=1) | 5.24 | 4.53× | 6.48 | 3.05 |
| ($\mathcal{N}$=6, $\mathcal{O}$=1) | 4.76 | 4.98× | 7.62 | 3.59 |
| ($\mathcal{N}$=7, $\mathcal{O}$=1) | 3.82 | 6.22× | 8.76 | 4.29 |
| ($\mathcal{N}$=3, $\mathcal{O}$=2) | 8.56 | 2.77× | 4.69 | 2.44 |
| ($\mathcal{N}$=4, $\mathcal{O}$=2) | 6.66 | 3.56× | 5.21 | 2.51 |
| ($\mathcal{N}$=5, $\mathcal{O}$=2) | 5.24 | 4.53× | 5.87 | 2.79 |
| ($\mathcal{N}$=6, $\mathcal{O}$=2) | 4.76 | 4.98× | 6.65 | 3.18 |
| ($\mathcal{N}$=7, $\mathcal{O}$=2) | 3.82 | 6.22× | 7.26 | 3.66 |
| ($\mathcal{N}$=3, $\mathcal{O}$=3) | 8.56 | 2.77× | 4.69 | 2.34 |
| ($\mathcal{N}$=4, $\mathcal{O}$=3) | 6.66 | 3.56× | 5.21 | 2.53 |
| ($\mathcal{N}$=5, $\mathcal{O}$=3) | 5.24 | 4.53× | 5.82 | 2.78 |
| ($\mathcal{N}$=6, $\mathcal{O}$=3) | 4.76 | 4.98× | 6.50 | 3.15 |
| ($\mathcal{N}$=7, $\mathcal{O}$=3) | 3.82 | 6.22× | 7.08 | 3.63 |
| ($\mathcal{N}$=3, $\mathcal{O}$=4) | 8.56 | 2.77× | 4.69 | 2.35 |
| ($\mathcal{N}$=4, $\mathcal{O}$=4) | 6.66 | 3.56× | 5.19 | 2.52 |
| ($\mathcal{N}$=5, $\mathcal{O}$=4) | 5.24 | 4.53× | 5.82 | 2.78 |
| ($\mathcal{N}$=6, $\mathcal{O}$=4) | 4.76 | 4.98× | 6.50 | 3.15 |
| ($\mathcal{N}$=7, $\mathcal{O}$=4) | 3.82 | 6.22× | 7.07 | 3.63 |

## 8. Anonymous Page for Video Presentation

To further showcase the advantages of TaylorCache in video generation, we have created an anonymous GitHub page.

For a more detailed demonstration, please visit https://taylorseer.github.io/TaylorSeer/. Additionally, the videos are also available in the Supplementary Material.

## 9. Supplementary Visualization Examples

To further illustrate the qualitative improvements of our method, we present visualization examples on FLUX and HunyuanVideo. These results showcase the superior fidelity and consistency of our method in generating high-quality outputs across diverse scenarios.

## 10. Supplementary Visualization of Feature Trajectories in Diffusion Models

In this section, we provide additional visualizations of feature trajectories and their derivatives in diffusion models. These results further illustrate the stability and predictability of feature dynamics across different timesteps, supporting our findings in the main text. The PCA projections of features (0th-order) and their derivatives (1st to 4th-order) demonstrate consistent patterns, highlighting the potential for efficient feature prediction in diffusion models.



Figure 8. **PCA projections of features in diffusion models.** The features at different timesteps form stable trajectories, demonstrating the predictability of feature evolution over time.
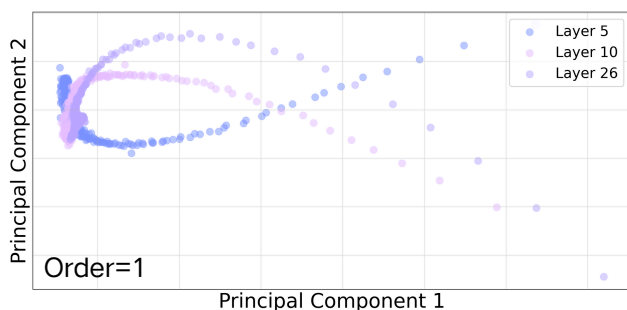


Figure 9. **PCA projections of first-order feature derivatives.** The first-order derivatives exhibit consistent patterns, further supporting the predictability of feature dynamics in diffusion models.
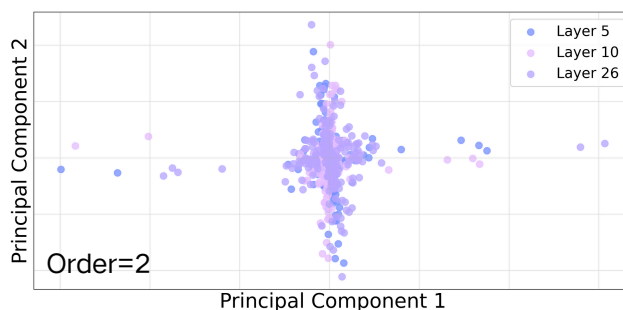


Figure 10. **PCA projections of second-order feature derivatives.** The second-order derivatives reveal higher-order dynamics, highlighting the smoothness of feature transitions.
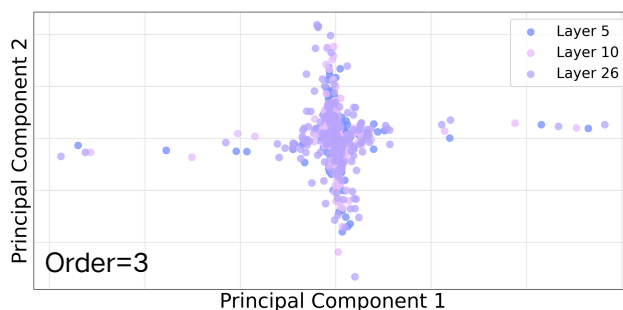


Figure 11. **PCA projections of third-order feature derivatives.** The third-order derivatives capture more complex temporal patterns, indicating the richness of feature dynamics.
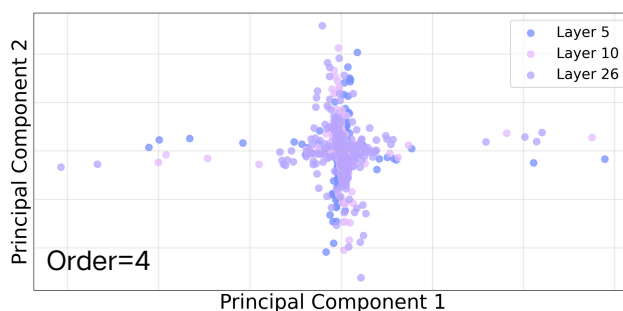


Figure 12. **PCA projections of fourth-order feature derivatives.** The fourth-order derivatives provide insights into fine-grained temporal variations, further validating the predictability of feature evolution.
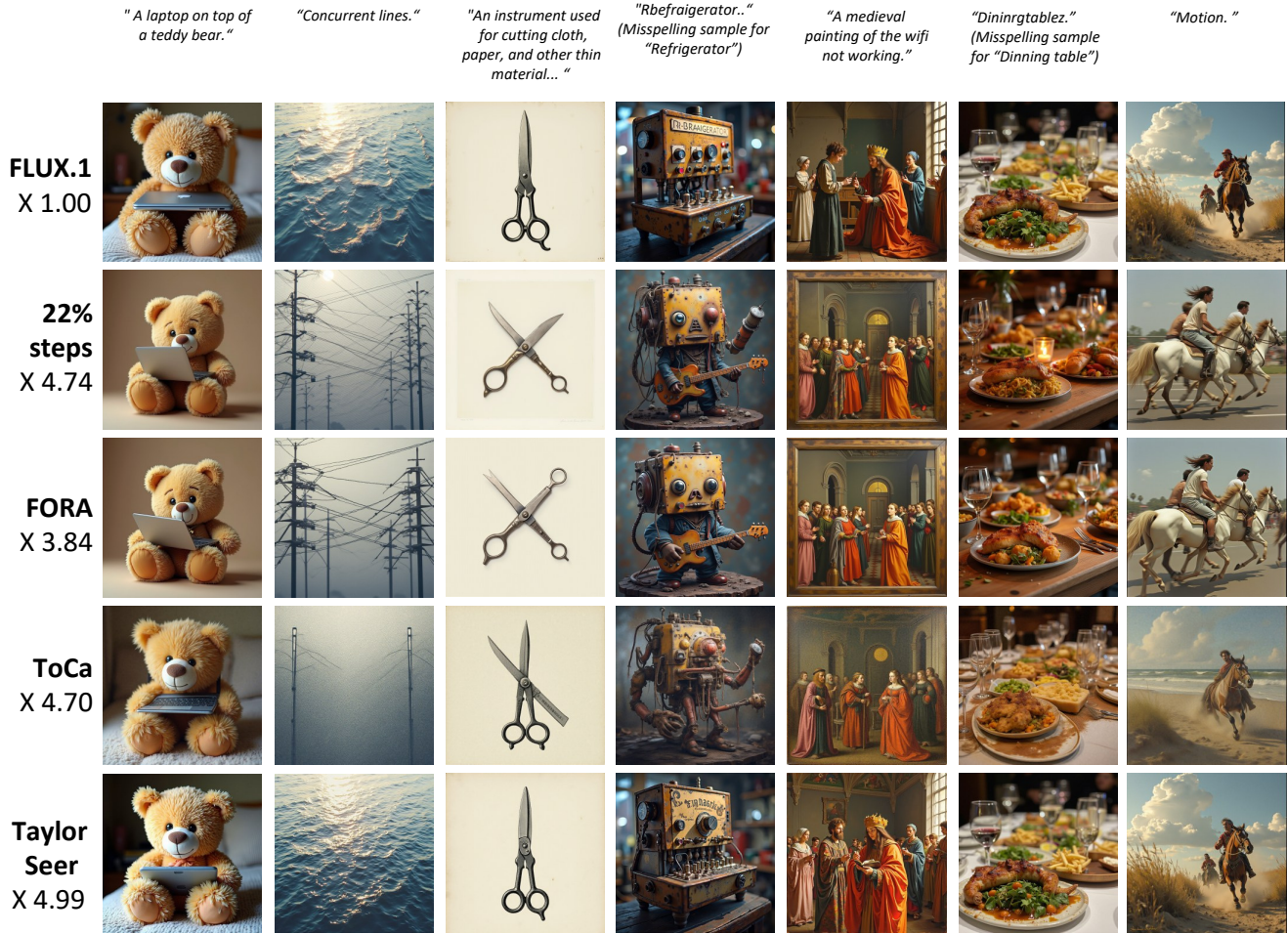
Figure 13. Visualization results for different acceleration methods on FLUX.1-dev.
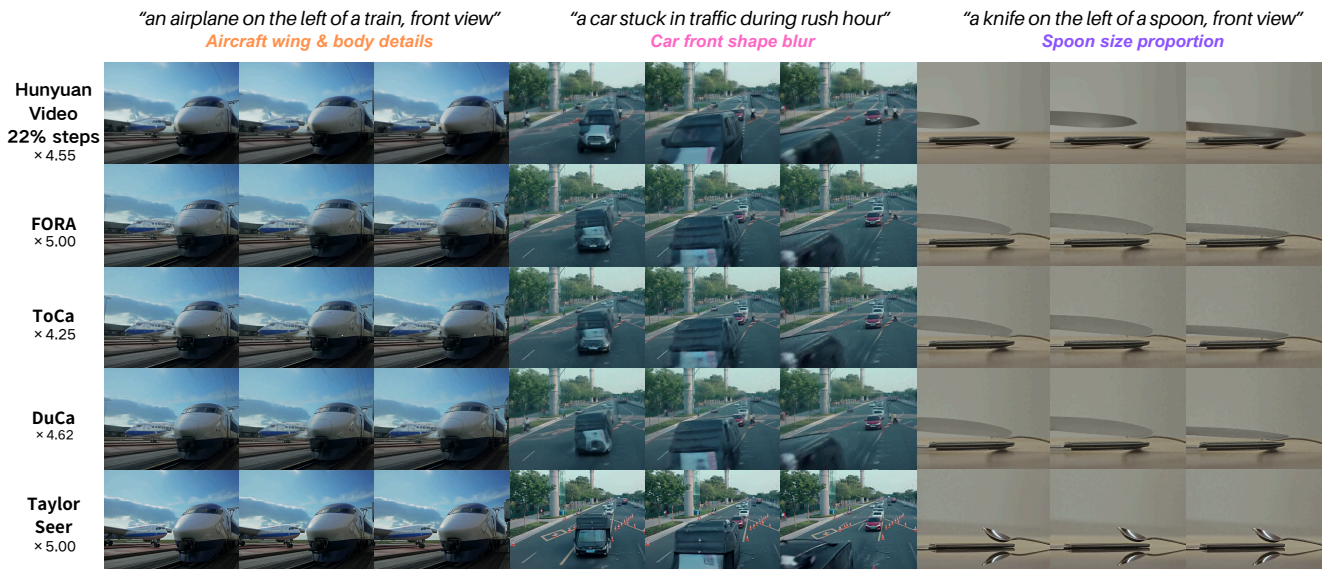


Figure 14. Visualization results for different acceleration methods on HunyuanVideo.