

DDiT: Dynamic Patch Scheduling for Efficient Diffusion Transformers

Anonymous CVPR submission

Paper ID ****

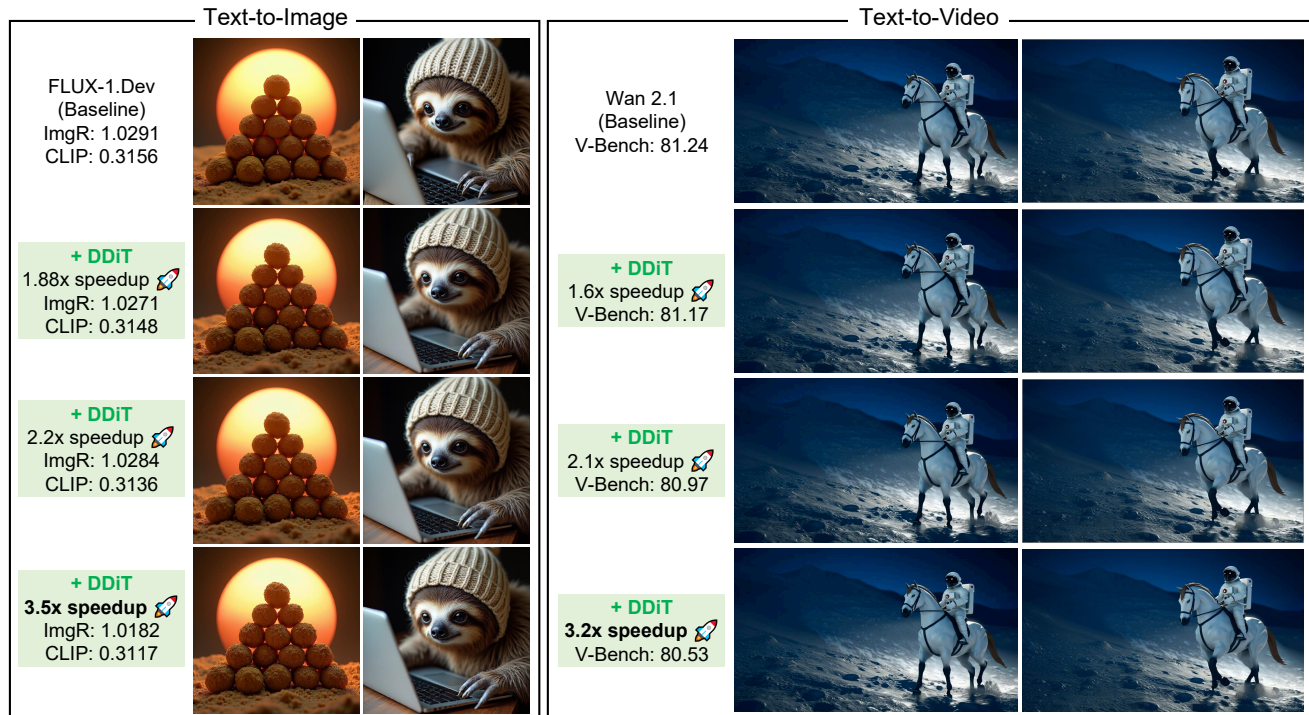


Figure 1. **DDiT dynamically selects the optimal patch size at each denoising step at inference yielding significant computational gains at no loss of perceptual quality.** Results are shown for FLUX-1.Dev [10] for text-to-image and Wan-2.1 [27] for text-to-video generation. The top panel denotes the baseline (original model), while the remaining panels illustrate outputs from DDiT at different acceleration rates. ImageReward [30], CLIP [20], and VBench [8] scores are reported (higher is better).

Abstract

001 Diffusion Transformers (DiTs) have achieved state-of-the-
 002 art performance in image and video generation, but their
 003 success comes at the cost of heavy computation. This ineffi-
 004 ciency is largely due to the fixed tokenization process, which
 005 uses constant-sized patches throughout the entire denoising
 006 phase, regardless of the content’s complexity. We propose
 007 dynamic tokenization, an efficient **test-time strategy** that
 008 varies patch sizes based on content complexity and the de-
 009 noising timestep. Our key insight is that early timesteps only
 010 require coarser patches to model global structure, while
 011 later iterations demand finer (smaller-sized) patches to re-
 012 fine local details. During inference, our method dynami-
 013 cally reallocates patch sizes across denoising steps for im-

age and video generation and substantially reduces cost 014
 while preserving perceptual generation quality. Extensive 015
 experiments demonstrate the effectiveness of our approach: 016
 it achieves up to 3.52× and 3.2× speedup on FLUX-1.Dev 017
 and Wan 2.1, respectively, without compromising the gen- 018
 eration quality and prompt adherence. 019

1. Introduction 020

Diffusion transformers (DiTs) [3, 19] have emerged as 021
 a dominant framework for content generation, producing 022
 high-quality and photorealistic results in both image and 023
 video synthesis. However, this impressive performance 024
 comes with substantial computational cost, significantly 025
 limiting the usage of these models in practice. Such high 026

027 computational demands of generative models have cata-
 028 pulted the development of more efficient generation meth-
 029 ods. Existing research has broadly focused on acceleration
 030 techniques such as feature caching [15, 32], feature prun-
 031 ing [1, 4], vector quantization [23, 24], and model distilla-
 032 tion [12, 22].

033 Although these approaches show promise, they suffer
 034 from two key limitations. First, many methods [5, 7, 15, 29]
 035 typically employ a **hard, static reduction strategy**, such as
 036 removing a fixed amount of weights, operations, or tokens.
 037 Such static approach can lead to significant quality degra-
 038 dation, as computations critical to a specific output might
 039 be permanently discarded [2, 29]. Second, most existing
 040 methods [15, 16, 29] apply a rigid, one-size-fits-all strategy,
 041 that is agnostic to the input. This is problematic, as dif-
 042 ferent prompts require varying levels of computational de-
 043 tail [17, 28]. A simple prompt like “a blue sky” should not
 044 require the same amount of computational resources com-
 045 pared to a prompt “a scene crowded with many zebras.” The
 046 rigidity in all existing solutions prevents us from dynam-
 047 ically allocating resources where they are needed most.

048 In this work, we address the rigid, one-size-fits-all com-
 049 putation of existing methods. Our approach is based on a
 050 key observation: the visual content generated by a diffusion
 051 model evolves at varied levels of detail. Some denoising
 052 timesteps establish coarse scene structure, while others re-
 053 fine fine-grained visual details. Recent studies [9, 11, 25,
 054 26] show that features generated at different timesteps of
 055 the denoising process encode different information, thus se-
 056 lecting the right timestep of diffusion features is important
 057 for successful downstream tasks such as classification [11]
 058 and visual reasoning [9]. Furthermore, [18, 28] note that
 059 this information can also be used for image editing, catering
 060 to different levels of detail in an image [17, 18, 28], and for
 061 generating more prompt-aligned images by injecting differ-
 062 ent levels of prompt information at different timesteps [21].

063 This leads us to a critical question: **should every denois-**
 064 **ing step process the latent at the same granularity?** Or,
 065 could some steps operate on a coarser latent, thereby yield-
 066 ing computational benefits, while others use a finer latent to
 067 preserve detail? Thus, unlike prior works [1, 4] which ap-
 068 proach efficiency by discarding weights or operations, we
 069 *dynamically allocate* it. Specifically, at every timestep, we
 070 adjust the patch size of the latent and adaptively use larger
 071 patches (coarser granularity) when less detail is required
 072 and smaller patches (finer granularity) when high fidelity
 073 is needed.

074 This, however, raises a new question: **how do we deter-**
 075 **mine the optimal patch size at any given timestep and for**
 076 **any given prompt?** For this, we measure the *rate of change*
 077 *of the latent manifold* over time. We hypothesize that this
 078 rate correlates with the level of detail being generated. If
 079 the underlying latent evolves slowly within a short timestep

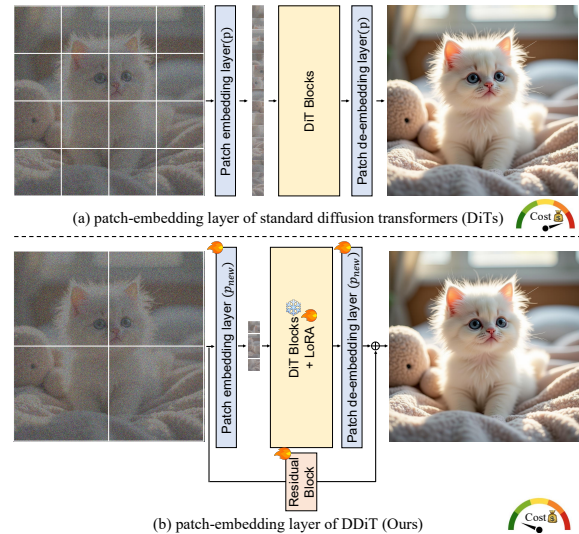


Figure 2. **Revised patch-embedding layer to support patches of varied resolutions.** We modify the standard patch-embedding layer, designed for a fixed patch size p , to additionally support patch sizes p_{new} .

080 window, we posit that coarse-grained details are being gener-
 081 ated. Consequently, we divide the latent into coarser
 082 patches and process them, saving computational resources.
 083 Conversely, if the underlying latent evolves rapidly, we infer
 084 that fine-grained details are being generated and fall back to
 085 using finer-grained latent patches.

086 Thus, this dynamic strategy tailors the computation load
 087 to each timestep and each prompt, allocating more re-
 088 sources when needed and conserving them where possible.
 089 Ultimately, our approach gives us explicit control over the
 090 computational budget while generating the highest possible
 091 quality content given the computational budget (Fig. 1). In
 092 summary, our contributions are:

- 093 • We introduce a simple, intuitive, and low-cost strategy to
 094 **dynamically vary latent’s granularity** in diffusion mod-
 095 els, that requires minimal architectural changes.
- 096 • We propose a test-time **Dynamic Patch Scheduler** that
 097 automatically determines the optimal patch size at each
 098 timestep, adapting the computational load based on gener-
 099 ation complexity and the input prompt.
- 100 • We demonstrate through extensive experiments that our
 101 approach generalizes across both image and video dif-
 102 fusion transformers and achieves significant speedups –
 103 up to $3.52\times$ on FLUX-1.Dev [10] and $3.2\times$ speedups on
 104 Wan 2.1 [27], while maintaining high perceptual quality,
 105 photo-realism, and prompt alignment.

106 2. Approach

107 2.1. Dynamic Patching and Tokenization

108 We aim to modify a pre-trained DiT to seamlessly oper-
 109 ate under different patch sizes with minimal architectural
 110 modifications (Fig. 2). To this end, we adapt the patch em-
 111 bedding layer, originally operating on patch size p , to also

112 handle new patch sizes p_{new} and allow input latents of vary-
113 ing spatial resolutions. We define p_{new} a positive integer
114 multiple of p , *i.e.*, $\{p, 2p, 4p, \dots\}$.

115 **Modifications to the patch embedding layer.** To gener-
116 alize DiT to p_{new} , we introduce patch-specific embedding
117 layers for each patch size we wish to support. Let C de-
118 notes the number of latent channels and d represents the
119 embedding dimension. Let $\mathbf{w}_{p_{\text{new}}}^{\text{emb}} \in \mathbb{R}^{p_{\text{new}} \times p_{\text{new}} \times C \times d}$ and
120 $\mathbf{b}_{p_{\text{new}}}^{\text{emb}} \in \mathbb{R}^d$ denote the weight matrix and bias vector of the
121 patch embedding layer corresponding to p_{new} . Each patch
122 of size p_{new} is linearly projected into an embedding of di-
123 mension d using this newly added embedding layer. This
124 results in a total of $N_{p_{\text{new}}} = \frac{HW}{p_{\text{new}}^2}$ patches. Since $N_{p_{\text{new}}}$
125 is smaller than N by a factor of $(p_{\text{new}}/p)^2$, DiT now processes
126 fewer patches and yields significant computational gains.

127 To minimize the training cost, we retain the base model
128 originally trained on the latent patch size p and introduce a
129 Low-Rank Adaptation (LoRA) branch [6] into each trans-
130 former block in DiT. This LoRA branch serves as an *adapt-*
131 *ive pathway* and enables the model to process patches of
132 different sizes. Additionally, as shown in Fig. 2, we add a
133 residual connection from *before* the patch embedding layer
134 to *after* the patch de-embedding block. This helps strike a
135 balance between the base latent manifold and the new man-
136 ifold being learnt by LoRA for p_{new} .

137 Finally, to distill the knowledge from the frozen base
138 model to the LoRA-augmented model, we fine-tune the
139 LoRA branch with a distillation loss. Let ϵ_{θ_L} and ϵ_{θ_T}
140 denote the predicted noise from the LoRA-fine-tuned and
141 frozen base models respectively. The distillation loss is:

$$142 \quad \mathcal{L} = \|\epsilon_{\theta_L}(\mathbf{z}_t^{p_{\text{new}}}, t) - \epsilon_{\theta_T}(\mathbf{z}_t^p, t)\|_2^2. \quad (1)$$

143 These minor architectural tweaks allow us to dynami-
144 cally support larger patch sizes while maintaining the base
145 model’s perceptual output quality.

146 2.2. Dynamic Patch Scheduling

147 Now that we enabled processing of multiple size input
148 patches, how do we learn *when* to adapt to a larger patch-
149 size (*i.e.*, coarser token) and when to switch back to a
150 smaller patch-size (*i.e.*, fine-grained token)? To this end, we
151 introduce a dynamic patch scheduling mechanism to deter-
152 mine the appropriate patch size at each diffusion timestep.
153 Since different timesteps correspond to different levels of
154 generative detail [17, 18, 21, 25, 26, 28], selecting the
155 proper patch size at each stage is crucial for maintaining
156 both efficiency and quality. We design a **training-free**
157 **dynamic scheduler** that adaptively selects the patch size
158 based on the *rate of evolution of the latent representations*
159 within a window of timesteps.

160 **Latent evolution estimation.** We employ finite-difference
161 approximations of increasing order to quantify how latent
162 representations evolve during the denoising process. Let \mathbf{z}_t

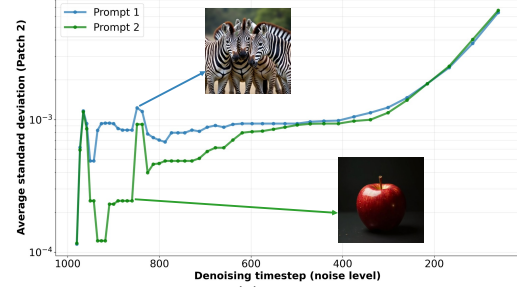


Figure 3. **Visualization of $\sigma_{t-1}^{2p, (\rho)}$ for two prompts (log scale).** *Prompt 1:* “Several zebras are standing together behind a fence.” *Prompt 2:* “A simple red apple on a black background.” Prompts requiring different levels of spatial granularity exhibit distinct $\sigma_{t-1}^{2p, (\rho)}$ patterns across timesteps. For the fine-grained zebra pattern, $\sigma_{t-1}^{2p, (\rho)}$ remains higher, indicating higher detail sensitivity, whereas for the simpler apple scene, $\sigma_{t-1}^{2p, (\rho)}$ is lower, thus we can seamlessly use larger patch sizes during generation.

163 denote the latent at timestep t . The first-order finite differ-
164 ence captures the displacement of latent features between
165 consecutive timesteps:

$$166 \quad \Delta \mathbf{z}_t = \mathbf{z}_t - \mathbf{z}_{t+1}. \quad (2)$$

167 Similarly, the second-order difference describes the *rate* of
168 change of this displacement, representing the local velocity
169 of the denoising trajectory, defined by

$$170 \quad \Delta^{(2)} \mathbf{z}_{t-1} = \Delta \mathbf{z}_{t-1} - \Delta \mathbf{z}_t. \quad (3)$$

171 Finally, the third-order finite difference quantifies the varia-
172 tion in this velocity. This can be interpreted as a measure of
173 *acceleration* of a latent’s evolution during denoising within
174 a short temporal window.

$$175 \quad \Delta^{(3)} \mathbf{z}_{t-1} = \Delta^{(2)} \mathbf{z}_{t-1} - \Delta^{(2)} \mathbf{z}_t = 2 \left(\frac{\Delta \mathbf{z}_{t-1} + \Delta \mathbf{z}_{t+1}}{2} - \Delta \mathbf{z}_t \right), \quad (4)$$

176 We hypothesize that if the acceleration is slow at a given
177 timestep, there is a relatively minor difference in the under-
178 lying latent manifold in the local temporal window. On the
179 other hand, **a high acceleration value suggests a larger
180 difference in the structure of the underlying manifold.**
181 We use this measure as a proxy to identify transition points
182 where the generative process intrinsically shifts between
183 generating coarse to fine structures or vice versa.

184 Empirically, we find that the third-order difference cap-
185 tures this variation more effectively and remains more stable,
186 while the first- and second-order differences fail to do
187 so, likely because they capture relatively short-term tempo-
188 ral changes. This observation is consistent with [31], which
189 shows that the difference between neighboring noise predic-
190 tions is explicitly related to the third-order finite difference.

191 **Spatial variance estimation.** We use latent \mathbf{z}_t in the above
192 formulation for simplicity, but in practice, \mathbf{z}_t is always di-
193 vided into patches of size $p_{\text{new}} \times p_{\text{new}}$ (Sec. 2.1). Our final
194 task now is to select the right patch size at each latent man-
195 ifold. This requires quantifying and aggregating the accel-
196 eration at which the latent patches evolve. Thus, we divide

Table 1. **Quantitative comparison of text-to-image generation performance with state-of-the-art methods** on COCO, DrawBench, and PartiPrompts. If not specified, all results are reported using 50 inference steps by default. Each color (**Yellow** , **Blue**) indicates methods operating at similar inference speeds. As highlighted in **Blue** , our method achieves the best overall image quality, evidenced by the lowest FID scores, strong prompt alignment (CLIP and ImageReward), and high perceptual similarity (SSIM and LPIPS). **Bold**: best. Underline: second-best.

Model	Speed↓ (secs/image)	COCO		DrawBench				PartiPrompts	
		FID↓	CLIP↑	CLIP↑	ImgR↑	SSIM↑	LPIPS↓	CLIP↑	ImgR↑
FLUX-1.Dev (50 steps)	12.0	33.07	0.314	0.3156	1.0291	–	–	0.3197	1.192
FLUX-1.Dev (28 steps)	6.72	33.35	0.312	0.3140	1.0107	0.739 ± 0.155	0.175 ± 0.126	0.3118	1.115
FLUX-1.Dev (15 steps)	3.6	34.02	0.311	0.3121	0.9865	0.591 ± 0.155	0.298 ± 0.128	0.3072	0.9613
TeaCache ($\delta = 0.6$)	6.0	34.95	0.303	0.3071	0.9968	0.654 ± 0.145	0.241 ± 0.114	0.3018	0.9701
TeaCache ($\delta = 0.8$)	5.33	35.57	0.303	0.3075	0.9780	0.612 ± 0.146	0.279 ± 0.114	0.2991	0.9699
TaylorSeer ($N = 3, O = 2$)	6.0	34.74	0.303	0.3085	0.9721	0.632 ± 0.173	0.271 ± 0.149	0.3142	0.9813
TaylorSeer ($N = 6, O = 1$)	3.5	35.02	0.302	0.3040	0.9535	0.583 ± 0.130	0.284 ± 0.106	0.3004	0.9714
DDiT	5.5	33.42	0.317	0.3136	1.0284	0.635 ± 0.183	0.264 ± 0.190	0.3192	1.189
DDiT + Teacache ($\delta = 0.4$)	3.4	<u>33.60</u>	<u>0.315</u>	<u>0.3117</u>	<u>1.0182</u>	0.592 ± 0.118	0.267 ± 0.120	<u>0.3172</u>	<u>1.062</u>

z_{t-1} into patches of size $p_i \times p_i$, where $p_i \in p_{\text{new}}$. Then, we compute the **standard deviation** $\sigma_{t-1}^{p_i}$ of the acceleration (defined in Eqn. 4) within each patch. We hypothesize that if the per-(latent) pixel standard deviation within a latent patch is high, then the denoising process is focusing on generating finer-grained details. On the other hand, if the standard deviation is low, then the underlying evolving latent is smooth. As shown in Fig. 3, prompts with different levels of granularity exhibit distinct variance profiles across timesteps.

Given $\sigma_{t-1}^{p_i}$, our goal is to determine the appropriate patch size at each timestep. A straightforward way to do this is to aggregate $\sigma_{t-1}^{p_i}$ by taking their mean across patches, which provides a simple measure of the overall latent variation at that timestep. However, this fails to effectively capture the generative dynamics occurring at that timestep. For example, when generating an image containing both a uniform white background and a highly textured region, averaging might smoothen the higher standard deviation values, leading to the scheduler choosing larger patches and thus overlooking fine details in the textured area. To better capture such spatial heterogeneity in the underlying latent manifold, we instead take the ρ -th percentile of the per-patch variances, denoted as $\sigma_{t-1}^{p_i,(\rho)}$. This percentile-based aggregation allows us to capture meaningful information across patches without averaging out important signals, while also avoiding bias toward a few high-variance outliers.

Concretely, we compare $\sigma_{t-1}^{p_i,(\rho)}$ against a predefined variance threshold τ . For each timestep, we select the largest patch size whose corresponding variance is below the threshold (τ). If no such patch satisfies this condition, it defaults to the smallest patch size, which is 1. We formulate this **patch size scheduling** as follows:

$$p_t = \begin{cases} \max(p_i), & \text{if } \sigma_{t-1}^{p_i,(\rho)} < \tau, \\ 1, & \text{otherwise.} \end{cases} \quad (5)$$

Controlling τ gives us explicit control over the speed: if users prefer faster generation, a higher τ can be selected;

otherwise, a smaller τ can be used for higher quality with less speed gain. We select τ and ρ empirically and balance generation stability and visual quality.

3. Experiments

3.1. Text-to-Image Generation

We evaluate our method on text-to-image generation using FLUX-1.dev [10] and compare against state-of-the-art acceleration methods [13, 14]. Our approach achieves comparable FID and CLIP scores to the base model while providing a 2.18 \times speedup, consistently outperforming prior methods under similar computational budgets. Furthermore, combining our method with TeaCache [13] yields up to a 3.52 \times speedup, demonstrating that dynamic patch scheduling improves efficiency while preserving perceptual quality.

3.2. Text-to-Video Generation

We further evaluate the effectiveness of our method in the text-to-video (T2V) generation setting and compare it with the base model [27]. As shown in Fig. 1, our approach significantly reduces inference time while maintaining competitive video quality, as reflected by the VBench score [8].

4. Conclusion

We present an intuitive and highly-computationally efficient method, *DDiT*, to adapt diffusion transformers to patches of different sizes during denoising while maintaining visual quality. *DDiT* demonstrates a critical insight: not all timesteps require the underlying latent space to be equally fine-grained. Building on this insight, we dynamically select the optimal patch size at every timestep and achieve significant computational gains, with no loss in perceptual visual quality. Our approach requires just adding a simple plug-and-play LoRA adapter to make the patch-embedding (and de-embedding) blocks amenable to varied input patch sizes. This minimal architectural tweak allows any DiT-based model to benefit from fast inference.

268

References

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

- [1] Daniel Bolya and Judy Hoffman. Token merging for fast stable diffusion. In *CVPR*, 2023. 2
- [2] Juncan Deng, Shuaiting Li, Zeyu Wang, Hong Gu, Kedong Xu, and Kejie Huang. Vq4dit: Efficient post-training vector quantization for diffusion transformers. In *AAAI*, 2025. 2
- [3] 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 *ICLR*, 2024. 1
- [4] Gongfan Fang, Xinyin Ma, and Xinchao Wang. Structural pruning for diffusion models. In *NeurIPS*, 2023. 2
- [5] Gongfan Fang, Kunjun Li, Xinyin Ma, and Xinchao Wang. Tinyfusion: Diffusion transformers learned shallow. In *CVPR*, 2025. 2
- [6] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *ICLR*, 2022. 3
- [7] Taihang Hu, Linxuan Li, Joost van de Weijer, Hongcheng Gao, Fahad Shahbaz Khan, Jian Yang, Ming-Ming Cheng, Kai Wang, and Yaxing Wang. Token merging for training-free semantic binding in text-to-image synthesis. In *NeurIPS*, 2024. 2
- [8] 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 *CVPR*, 2024. 1, 4
- [9] Dahye Kim, Xavier Thomas, and Deepti Ghadiyaram. Revelio: Interpreting and leveraging semantic information in diffusion models. In *ICCV*, 2025. 2
- [10] Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024. 1, 2, 4
- [11] Alexander C Li, Mihir Prabhudesai, Shivam Duggal, Ellis Brown, and Deepak Pathak. Your diffusion model is secretly a zero-shot classifier. In *ICCV*, 2023. 2
- [12] 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. In *NeurIPS*, 2023. 2
- [13] 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. In *CVPR*, 2025. 4
- [14] Jiacheng Liu, Chang Zou, Yuanhuiyi Lyu, Junjie Chen, and Linfeng Zhang. From reusing to forecasting: Accelerating diffusion models with taylorseers. *arXiv preprint arXiv:2503.06923*, 2025. 4
- [15] Xinyin Ma, Gongfan Fang, Michael Bi Mi, and Xinchao Wang. Learning-to-cache: Accelerating diffusion transformer via layer caching. In *NeurIPS*, 2024. 2
- [16] Xinyin Ma, Gongfan Fang, and Xinchao Wang. Deepcache: Accelerating diffusion models for free. In *CVPR*, 2024. 2
- [17] Shweta Mahajan, Tanzila Rahman, Kwang Moo Yi, and Leonid Sigal. Prompting hard or hardly prompting: Prompt inversion for text-to-image diffusion models. In *CVPR*, 2024. 2, 3
- [18] Or Patashnik, Daniel Garibi, Idan Azuri, Hadar Averbuch-Elor, and Daniel Cohen-Or. Localizing object-level shape variations with text-to-image diffusion models. In *ICCV*, 2023. 2, 3
- [19] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *ICCV*, 2023. 1
- [20] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 1
- [21] Ketan Suhaas Saichandran, Xavier Thomas, Prakhara Kaushik, and Deepti Ghadiyaram. Progressive prompt detailing for improved alignment in text-to-image generative models. *arXiv preprint arXiv:2503.17794*, 2025. 2, 3
- [22] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022. 2
- [23] Yuzhang Shang, Zhihang Yuan, Bin Xie, Bingzhe Wu, and Yan Yan. Post-training quantization on diffusion models. In *CVPR*, 2023. 2
- [24] Junhyuk So, Jungwon Lee, Daehyun Ahn, Hyungjun Kim, and Eunhyeok Park. Temporal dynamic quantization for diffusion models. In *NeurIPS*, 2023. 2
- [25] Nick Stracke, Stefan Andreas Baumann, Kolja Bauer, Frank Fundel, and Björn Ommer. Cleandiff: Diffusion features without noise. In *CVPR*, 2025. 2, 3
- [26] Luming Tang, Menglin Jia, Qianqian Wang, Cheng Perng Phoo, and Bharath Hariharan. Emergent correspondence from image diffusion. In *NeurIPS*, 2023. 2, 3
- [27] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025. 1, 2, 4
- [28] Luozhou Wang, Shuai Yang, Shu Liu, and Ying-cong Chen. Not all steps are created equal: Selective diffusion distillation for image manipulation. In *ICCV*, 2023. 2, 3
- [29] Zhendong Wang, Yifan Jiang, Huangjie Zheng, Peihao Wang, Pengcheng He, Zhangyang Wang, Weizhu Chen, Mingyuan Zhou, et al. Patch diffusion: Faster and more data-efficient training of diffusion models. In *NeurIPS*, 2023. 2
- [30] 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. In *NeurIPS*, 2023. 1
- [31] Hancheng Ye, Jiakang Yuan, Renqiu Xia, Xiangchao Yan, Tao Chen, Junchi Yan, Botian Shi, and Bo Zhang. Training-free adaptive diffusion with bounded difference approximation strategy. In *NeurIPS*, 2024. 3
- [32] Hui Zhang, Tingwei Gao, Jie Shao, and Zuxuan Wu. Blockdance: Reuse structurally similar spatio-temporal features to accelerate diffusion transformers. In *CVPR*, 2025. 2