

# CoReDiT: Spatial Coherence-Guided Token Pruning and Reconstruction for Efficient Diffusion Transformers

Anonymous CVPR submission

Paper ID 9

## Abstract

001 *Diffusion Transformers (DiTs) deliver remarkable im-*  
002 *age and video generation quality but incur high compu-*  
003 *tational cost, limiting scalability and on-device deploy-*  
004 *ment. We introduce CoReDiT, a structured token prun-*  
005 *ing framework for DiTs across vision tasks. CoReDiT*  
006 *uses a linear-time spatial coherence score to estimate lo-*  
007 *cal redundancy in the latent token lattice and skips high-*  
008 *coherence (redundant) tokens in self-attention. To main-*  
009 *tain a dense representation and avoid visual discontinuities,*  
010 *we reconstruct skipped attention outputs via coherence-*  
011 *guided aggregation of spatially neighboring retained to-*  
012 *kens. We further introduce a progressive, block-adaptive*  
013 *pruning schedule that increases pruning gradually and al-*  
014 *locates larger budgets to blocks and denoising steps with*  
015 *higher redundancy. Across state-of-the-art diffusion back-*  
016 *bones including PixArt- $\alpha$  and MagicDrive-V2, CoReDiT*  
017 *achieves up to 55% self-attention FLOPs reduction and*  
018 *inference speedups of 1.33 $\times$  on cloud GPUs and 1.72 $\times$*   
019 *on mobile NPUs, while maintaining high visual quality.*  
020 *Notably, CoReDiT also increases on-device memory head-*  
021 *room, enabling higher-resolution generation.*

## 022 1. Introduction

023 Diffusion models learn to synthesize data through an iter-  
024 ative denoising process, and have achieved state-of-the-art  
025 results across various generative applications, such as con-  
026 ditional image generation, image editing (inpainting and  
027 super-resolution) and video synthesis. Building on this  
028 success, Diffusion Transformers (DiTs) [21] replace con-  
029 volutional U-Nets with attention-based backbones that ef-  
030 fectively model long-range dependencies among tokenized  
031 patches; this architecture scale effectively with model size  
032 and facilitate the incorporation of conditioning from various  
033 modalities (e.g., class labels, text, and depth maps).

034 However, these benefits come with demanding compu-  
035 tational cost. Self-attention scales quadratically with

the number of tokens, making high-resolution images and 036  
multi-frame videos particularly expensive and often im- 037  
practical on compute- and memory-constrained mobile plat- 038  
forms. Empirically, a large fraction of tokens correspond 039  
to visually redundant or low-saliency regions (e.g., uniform 040  
backgrounds and smooth textures), suggesting substantial 041  
room for compute reduction. Prior work reduces inference 042  
cost by pruning or merging tokens using attention-based 043  
saliency [12, 27], feature similarity [2], or learned rout- 044  
ing/predictors [23, 34]. However, several challenges limit 045  
their applicability to diffusion inference: (1) localizing low- 046  
saliency tokens efficiently and effectively, (2) preserving vi- 047  
sual semantics for the tokens that do not participate in at- 048  
tention, and (3) determining a pruning schedule that adapts 049  
across transformer blocks and denoising timesteps. 050

**Contributions:** To address these challenges, we propose 051  
CoReDiT, a token pruning framework for accelerating *pre-* 052  
*trained* DiTs that targets self-attention while preserving out- 053  
put fidelity. CoReDiT consists of three components: (1) 054  
*Spatial coherence-based selector* (Sec. 3.1): We observe 055  
that redundant tokens in low-saliency regions are highly 056  
similar to their spatial neighbors. To exploit this, we parti- 057  
tion the token lattice into small, non-overlapping grids and 058  
estimate a *spatial coherence* score for each token, based on 059  
its feature similarity to tokens in the same grid. This score 060  
can be implemented with linear-time, hardware-friendly re- 061  
ductions (e.g., grid means and dot products), adding min- 062  
imal overhead while reliably selecting redundant tokens 063  
(with high coherence) to bypass attention. (2) *Coherence-* 064  
*based reconstruction* (Sec. 3.2): Skipping tokens can break 065  
local spatial continuity, leading to artifacts in the generated 066  
outputs. To preserve a dense token lattice for subsequent 067  
blocks, we reconstruct skipped tokens from retained neigh- 068  
bors using similarity-weighted aggregation within a local 069  
neighborhood. This content-aware interpolation (unlike ze- 070  
roing or naive forwarding) preserves visual semantics and 071  
mitigates artifacts while maintaining compatibility with the 072  
original DiT architecture. (3) *Progressive, block-adaptive* 073  
*pruning* (Sec. 3.3): Pruning tolerance varies across blocks 074  
and timesteps: different blocks carry varying semantic im- 075

076 portance, and late diffusion steps are typically less tolerant  
077 than early ones [34]. Hence, we fine-tune with a progressive  
078 schedule that *automatically allocates* the pruning budget to  
079 the most redundant blocks. At regular intervals during fine-  
080 tuning, we estimate a timestep-weighted redundancy score  
081 for each transformer block and increment the pruning ratio  
082 of the block with the highest redundancy; gradually increas-  
083 ing the pruning ratio enables fast recovery after each update  
084 and ensures stable adaptation.

085 Combining these components results in an efficient  
086 pruning pipeline that reduces attention computation while  
087 preserving high-fidelity generation. Across text-to-image  
088 and video generation, CoReDiT enables higher-resolution  
089 synthesis and stable conditional alignment, achieving up  
090 to 55% reduction in self-attention FLOPs and wall-clock  
091 speedups of  $1.33\times$  on cloud GPUs and  $1.72\times$  on mobile  
092 devices.

## 093 2. Related Work

094 **Diffusion Models.** Diffusion models have demonstrated  
095 impressive generative performance across image, video,  
096 and 3D domains. Early works commonly adopt U-Net  
097 backbones [11, 25], and latent diffusion [24] improves scal-  
098 ability by denoising in a compressed latent space. Recently,  
099 Diffusion Transformers (DiTs) [21] replace U-Net archi-  
100 tectures with stacks of transformer blocks, offering better  
101 scalability to high resolution and flexible conditioning, but  
102 incurring high inference cost from quadratic attention and  
103 multi-step denoising.

104 **Efficient Diffusion Transformers.** Prior studies have  
105 improved the efficiency of DiTs along several comple-  
106 mentary directions: (1) *Token reduction within a timestep.*  
107 Training-free token reduction often merges or downsam-  
108 ples tokens to reduce the effective sequence length, in-  
109 cluding token merging methods (e.g., ToMe [1]) and  
110 diffusion-specific variants (e.g., ToMA [18]) that aim to  
111 reduce merge overhead. Other training-free strategies re-  
112 duce attention cost by downsampling only key/value to-  
113 kens (e.g., ToDo [28]). Learning-based methods learn per-  
114 layer/timestep keep ratios or token routers for content gen-  
115 eration, including DiffCR [34] and dynamic token-density  
116 schemes such as FlexDiT [5]. A key challenge in DiTs  
117 (vs. ViTs) is that naively discarding tokens can introduce  
118 spatial discontinuities (texture breaks, boundary artifacts),  
119 motivating approaches that either preserve structure via  
120 merge/unmerge or explicitly reconstruct skipped content.  
121 (2) *Cross-timestep caching and feature reuse.* Caching-  
122 based methods exploit temporal redundancy across de-  
123 noising steps by reusing intermediate features with se-  
124 lective refresh [19, 31]. Recent DiT-specific works fur-  
125 ther cache at finer granularity (e.g., token-wise feature  
126 caching in ToCa [39]). These approaches primarily tar-  
127 get *cross-step redundancy* and are conceptually orthogo-

128 nal to within-step token reduction. (3) *Dynamic compu-*  
129 *tation and conditional routing.* Dynamic DiT architec-  
130 tures adjust compute across timesteps and/or spatial re-  
131 gions, e.g., DyDiT [38] uses timestep-wise dynamic width  
132 and spatial token skipping, while DiffCR [34] fine-tunes to-  
133 ken routers with differentiable compression ratios. Com-  
134 pared to token reduction, these methods introduce explicit  
135 routing modules and typically optimize a dynamic com-  
136 pute graph. (4) *Sparse/structured attention kernels.* An-  
137 other line accelerates DiTs by replacing full attention with  
138 sparse attention patterns and hardware-efficient kernels, in-  
139 cluding training-free profiling-based sparse attention (e.g.,  
140 Sparse VideoGen [32]) and trainable sparse attention (e.g.,  
141 VSA [35]). These methods change the attention operator  
142 itself rather than reducing the token set, and can be com-  
143plementary to token pruning/merging. (5) *Weight-pruning*  
144 *and architecture editing.* Beyond tokens, structured and un-  
145 structured pruning remove channels, heads, or even full DiT  
146 blocks, with brief recovery fine-tuning or lightweight cali-  
147 bration to retain quality [8]. Other work (e.g., grafting [4])  
148 edits the architectures of pretrained DiTs to explore more  
149 efficient backbones under small compute budgets; these  
150 approaches are naturally complementary to token sparsity  
151 methods.

152 **Positioning of CoReDiT.** CoReDiT is a *post-training*  
153 method that targets *within-timestep* redundancy in DiT la-  
154 tent grids: it performs *structured token skipping* guided  
155 by local spatial coherence, and *reconstructs* skipped tokens  
156 to preserve a full lattice for subsequent layers. This dif-  
157 fers from caching-based methods (e.g., FasterCache/ToCa)  
158 that exploit *cross-timestep reuse*, from dynamic routing ap-  
159 proaches (e.g., FlexDiT/DiffCR) that introduce learnable  
160 routing policies, and from token-merging methods (e.g.,  
161 ToMe) that explicitly merge/unmerge tokens. Compared  
162 to purely training-free schemes, CoReDiT uses lightweight  
163 fine-tuning to learn block/timestep-specific pruning sched-  
164 ules to improve hardware-relevant speed/quality trade-offs.

## 165 3. Proposed Approach: CoReDiT

166 Fig. 1 summarizes the workflow of CoReDiT: In each trans-  
167 former block, we partition the input tokens (patch embed-  
168 dings)  $X$  into a retained set  $X_R$  and a skipped set  $X_S$  ac-  
169 cording to a spatial coherence score that quantifies local re-  
170 dundancy of each token; only  $X_R$  participates in multi-head  
171 self-attention, producing  $Y_R = \text{MHSA}(X_R)$  (Sec. 3.1).  
172 For skipped tokens  $X_S$ , we synthesize their attention out-  
173 puts  $Y_S$  using nearby retained outputs  $Y_R$  with coherence-  
174 guided weights (Sec. 3.2). Importantly, *attention among re-*  
175 *tained tokens remains global*; locality is used only for ef-  
176 ficient selection and reconstruction, not to impose a local-  
177 attention pattern. To determine which block to prune more  
178 aggressively, we learn a progressive, block-adaptive prun-  
179 ing schedule during fine-tuning of a pretrained DiT, based

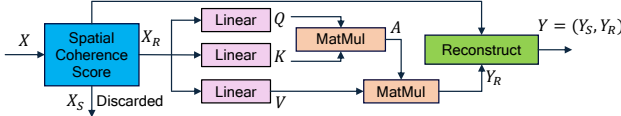


Figure 1. Workflow of CoReDiT. In each transformer block, input tokens  $X$  are partitioned into retained tokens  $X_R$  and skipped tokens  $X_S$  using the proposed spatial coherence score that captures *local redundancy* (Sec. 3.1). Only the retained tokens  $X_R$  participate in self-attention, reducing attention compute by skipping  $X_S$ . We then reconstruct the missing attention outputs  $Y_S$  from nearby retained outputs  $Y_R$  via coherence-guided aggregation (Sec. 3.2), preserving a dense token lattice for subsequent blocks. Finally, we learn an adaptive pruning-ratio schedule via coherence-guided progressive pruning during lightweight fine-tuning (Sec. 3.3).

180 on the estimated redundancy statistics (Sec. 3.3).

### 181 3.1. Token Selection

182 **Motivation.** A common token pruning criterion is at- 226  
 183 tention score-based importance [12, 27]: Let  $A =$  227  
 184  $\text{softmax}(\frac{QK^T}{\sqrt{d_k}})$  denote the attention matrix, where the 228  
 185 column-sum in  $A$  measures the total attention a token 229  
 186 receives. As Fig. 2a shows, pruning high-attention to- 230  
 187 kens severely damages visual semantics (salient regions), 231  
 188 whereas pruning low-attention tokens tends to be less harm- 232  
 189 ful (background regions). However, computing global at- 233  
 190 tention score  $A$  requires quadratic time and memory in to- 234  
 191 ken count, which substantially offsets pruning gains. Our 235  
 192 key observation is that low-saliency tokens typically exhibit 236  
 193 *strong local redundancy*: their embeddings are highly sim- 237  
 194 ilar to spatially nearby tokens. Motivated by this, we de- 238  
 195 sign an efficient selection score based on local similarity 239  
 196 that avoids forming the full  $N \times N$  attention matrix.

197 **Spatial Coherence (SC).** Given  $N$  tokens  $X = \{x_i\}_{i=1}^N$  240  
 198 reshaped into a 2D spatial lattice  $H \times W$  (given  $N = HW$ ), 241  
 199 we partition tokens into non-overlapping spatial grids of 242  
 200 size  $w \times w$  to estimate local (rather than global) similarity 243  
 201 (Fig. 2b). For a token  $x_i$ , define its spatial coherence (SC) 244  
 202 as the mean cosine similarity to other tokens in the same 245  
 203 grid:

$$204 \quad \text{SC}(x_i) := \frac{1}{|\mathcal{G}(i)|} \sum_{j \in \mathcal{G}(i)} \text{sim}(x_i, x_j) \quad (1)$$

205 where  $\mathcal{G}(i)$  indexes tokens in the same grid as token  $x_i$ . 250  
 206 Specifically, we adopt cosine similarity  $\text{sim}(x_i, x_j) =$  251  
 207  $\frac{x_i^T x_j}{\|x_i\|_2 \|x_j\|_2}$  due to its computational efficiency and strong 252  
 208 empirical performance in prior work [2]. Averaging within 253  
 209 each grid normalizes SC values across blocks with differ- 254  
 210 ent  $w$ , enabling comparisons in our block-adaptive schedule 255  
 211 (Sec. 3.3).

**Efficient computation and overhead.** Naively comput- 212  
 213 ing within-grid pairwise similarities costs  $O(w^2N)$  dot 214  
 215 products. Instead, let  $\hat{x}_i = \frac{x_i}{\|x_i\|_2}$  denote the normal- 216  
 217 ized token embedding, and  $\hat{g}_i = \frac{1}{|\mathcal{G}(i)|} \sum_{j \in \mathcal{G}(i)} \hat{x}_j$  denote 218  
 the mean normalized embeddings within the grid of token  $x_i$ . Then we can compute the spatial coherence score efficiently:

$$219 \quad \text{SC}(x_i) = \frac{1}{|\mathcal{G}(i)|} \sum_{j \in \mathcal{G}(i)} \hat{x}_i^T \hat{x}_j = \frac{\hat{x}_i^T}{|\mathcal{G}(i)|} \sum_{j \in \mathcal{G}(i)} \hat{x}_j = \hat{x}_i^T \hat{g}_i. \quad (2)$$

220 Essentially, the spatial coherence is an inner product be- 221  
 222 tween a token’s normalized embedding and the mean nor- 223  
 224 malized embedding of its grid. Thus, SC can be computed 225  
 with one reduction per grid plus one dot product per token, i.e.,  $O(N)$  time, which is linear in token count, regardless of grid size  $w$ .

**Selection rule.** High SC indicates a token is well- 226  
 227 explained by nearby patches (redundant), so we skip the 228  
 top- $K$  tokens by SC:

$$229 \quad X_S = \text{Top}_K(\{\text{SC}(x_i)\}_{i=1}^N), \quad X_R = X \setminus X_S \quad (3)$$

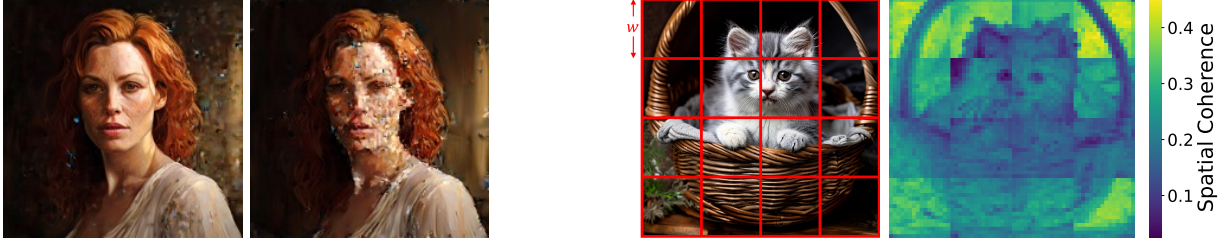
230 where  $K = r(l, s) \cdot N$  is determined by the pruning ratio 231  
 232  $r(l, s)$  for block  $l$  and timestep  $s$ . Empirically, as shown 233  
 234 in Fig. 2b, high-SC tokens (brighter colors) concentrate in 235  
 236 backgrounds, while boundaries/textures tend to have lower SC. Note that our proposed spatial coherence is not only used for token selection, but also utilized to update skipped tokens and guide our progressive pruning (Sec. 3.3).

### 237 3.2. Token Reconstruction

238 **Motivation.** Unlike ViTs for discriminative tasks (e.g., 239  
 240 classification and object detection) where dropping tokens 241  
 242 can be tolerated [15], generative DiTs require spatially co- 243  
 244 herent representations at every block and timestep. Naively 244  
 245 skipping tokens in MHSA can introduce discontinuities (in- 245  
 246 consistent textures, distorted boundaries), as in Fig. 3a. 246  
 247 We therefore reconstruct the missing attention outputs for 247  
 248 skipped tokens using nearby retained outputs, preserving 248  
 249 a dense lattice for subsequent blocks. Because our selec- 249  
 tor preferentially skips *redundant* (high-SC) tokens, these tokens are particularly well-suited to reconstruction from neighbors.

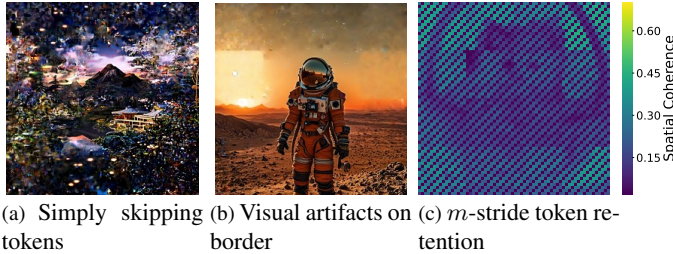
**Coherence-based reconstruction.** Let  $X_R$  be tokens re- 250  
 251 tained for attention and  $Y_R = \text{MHSA}(X_R)$ . For a skipped 252  
 253 token  $x_i \in X_S$ , a natural reconstruction is similarity- 254  
 255 weighted aggregation over nearby retained tokens:

$$254 \quad y_i = \frac{\sum_{j \in \mathcal{S}(i)} \text{sim}(x_i, x_j) y_j}{\sum_{k \in \mathcal{S}(i)} \text{sim}(x_i, x_k)} \quad (4)$$



(a) Token pruning with attention scores-based importance: removing tokens with the lowest scores (left) vs. the highest scores (right). (b) Example generated from PixArt- $\alpha$ -1024. Left: image with grid partition overlay. Right: computed SC scores over  $64 \times 64$  tokens.

Figure 2. Token selection motivation and visualization.



(a) Simply skipping tokens (b) Visual artifacts on border (c)  $m$ -stride token retention

Figure 3. Issues and proposed improvements. (a) Naively skipping tokens causes inconsistent textures, which can be mitigated by token reconstruction. (b) Using a fixed grid size during partitioning causes visual artifacts along grid borders; alternating grid sizes across consecutive blocks mitigates these artifacts. (c)  $m$ -stride token retention ensures at least one token is retained in each sub-grid to serve as a reconstruction anchor.

255 where  $\mathcal{S}(i)$  denotes retained tokens in a local neighborhood  
 256 of  $x_i$ . Intuitively, the transformation result of a token should  
 257 be more alike to that of a token with a higher similarity.  
 258 However, computing  $\text{sim}(x_i, x_j)$  per pair adds non-trivial  
 259 overhead.

260 **Approximate  $\text{sim}(x_i, x_j)$  with  $\text{SC}(x_j)$ .** Skipped tokens  
 261 are selected to have *high* coherence within their grid  
 262 (Sec. 3.1), meaning  $\hat{x}_i$  is close to the grid mean  $\hat{g}_i$ . For-  
 263 mally, define  $\epsilon_i = \|\hat{x}_i - \hat{g}_i\|_2$ ; high-SC implies small  $\epsilon_i$ .  
 264 Then for any  $j \in \mathcal{G}(i)$ :

$$\text{sim}(x_i, x_j) = \hat{x}_i^T \hat{x}_j = \hat{g}_i^T \hat{x}_j + (\hat{x}_i - \hat{g}_i)^T \hat{x}_j = \text{SC}(x_j) + O(\epsilon_i) \quad (5)$$

265 since  $\|\hat{x}_j\|_2 = 1$  bounds the residual term by  $\epsilon_i$ . Thus,  
 266 for the tokens we prune (high-SC, small  $\epsilon_i$ ), using  $\text{SC}(x_j)$   
 267 yields a close approximation to similarity-based weights  
 268 while avoiding pairwise similarity computation. We there-  
 269 fore approximate Eq. (4) as:  
 270

$$271 \quad y_i \approx \frac{\sum_{j \in \mathcal{S}(i)} \text{SC}(x_j) y_j}{\sum_{k \in \mathcal{S}(i)} \text{SC}(x_k)}. \quad (6)$$

**Sub-grids and complexity.** If  $\mathcal{S}(i)$  spans an entire grid,  
 Eq. (6) may over-smooth and reduce spatial variation. We  
 address this by subdividing each  $w \times w$  grid into smaller  
*sub-grids* of size  $w_s \times w_s$  and defining  $\mathcal{S}(i)$  as retained  
 tokens in the *same sub-grid* as  $x_i$ . This preserves local-  
 ity at boundaries while keeping computation linear. Re-  
 construction uses only local weighted sums over at most  
 $w_s^2$  tokens per skipped token; with small  $w_s$  (e.g., 3), the  
 overhead is modest and scales linearly with token count,  
 i.e.,  $O(N)$ . Notably, this locality is *only* for reconstruction;  
 MHSA among  $X_R$  is still global, so CoReDiT does not im-  
 pose a local attention pattern.

**Micro designs.** (1) Alternating grid size: Fixed grid  
 boundaries can introduce border artifacts due to limited  
 inter-grid mixing in reconstruction (Fig. 3b). We mitigate  
 this by alternating grid sizes across blocks (e.g., for  $64 \times 64$   
 tokens, alternate grid size  $16 \times 16$  and  $9 \times 9$ ), so boundaries  
 shift and information can propagate across grid borders over  
 consecutive blocks. This strategy enables cross-border in-  
 formation flow over consecutive blocks, eliminating visible  
 discontinuities. (2)  $m$ -stride token retention: In highly re-  
 dundant regions, an entire sub-grid could be skipped, leav-  
 ing no anchors for reconstruction. We therefore enforce re-  
 taining at least one token per  $m$  positions ( $m \leq w_s$ ) by  
 protecting a deterministic stride pattern. For block index  
 $l$ , we retain tokens at spatial coordinates  $[i, j]$  satisfying  
 $(i + j - l) \bmod m = 0$  by subtracting a large offset from  
 their SC so they are never selected into  $X_S$ .

### 3.3. Pruning Ratio Schedule

**Motivation.** Token redundancy varies across both de-  
 noising timesteps and transformer blocks. Fig. 4 depicts  
 the evolution of spatial coherence across denoising steps  
 and blocks. As shown, early steps display broadly ele-  
 vated coherence (coarse structure), mid steps accentuate  
 background redundancy, and late steps concentrate low-  
 coherence pockets around edges and details [5, 34]; besides,  
 within a timestep, different blocks exhibit different degrees  
 of locality [4]. Motivated by this, we propose adapting the

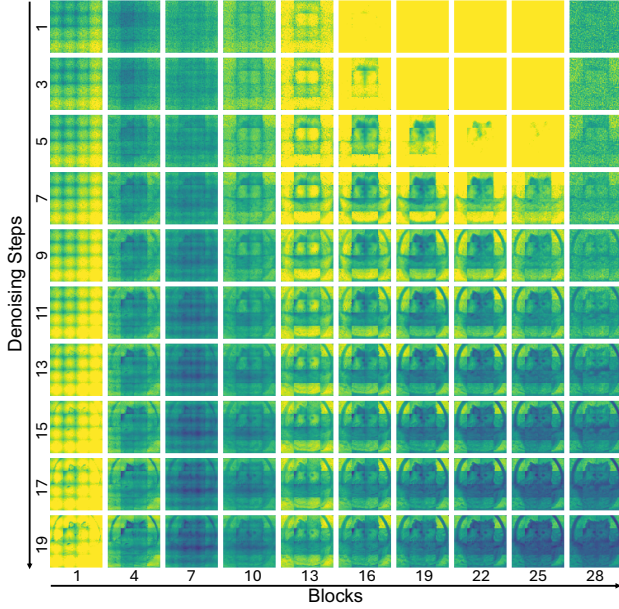


Figure 4. SC score evolution across transformer blocks and denoising timesteps (brighter = higher SC) for generating the image in Fig. 2b. (a) Early timestep in the first block: representations are noise-dominated and exhibit high local redundancy. (b) Early timestep in the last block: weaker spatial structure and reduced locality. (c) Late timestep in the first block: refined spatial details with clearer object-level coherence. (d) Late timestep in the last block: partial refinement with moderate spatial structure. These patterns indicate that redundancy varies with both timestep and block, motivating adaptive pruning ratios across depth and denoising phase.

310 pruning ratio across blocks based on the statistics observed  
 311 from real data. Specifically, we adopt progressive pruning  
 312 during fine-tuning: gradually increasing pruning bud-  
 313 gets while prioritizing blocks/timesteps with higher mea-  
 314 sured redundancy. This progressive design also conforms  
 315 to the empirical finding in [33]: Applying a large pruning  
 316 ratio abruptly to a pretrained model can be hard to recover,  
 317 while progressive pruning leads to better results at compa-  
 318 rable pruning levels.

319 **Coherence-based progressive pruning.** Let  $K(l, s) =$   
 320  $r(l, s) \cdot N$  denote the number of tokens pruned in block  $l$   
 321 at denoising step  $s$ . During fine-tuning, we increase  $K(l, s)$   
 322 in increments of  $\Delta k$  every  $T$  training iterations until reaching  
 323 a target budget. To decide which block to increase next, we  
 324 sort the SC scores in block  $l$  in descending order,  $\{\text{SC}_i^l\}_{i=1}^N$ .  
 325 We estimate the redundancy of the *next*  $\Delta k$  candidate to-  
 326 kens as:

$$327 \quad \Delta R_l = \sum_{i=K(l,s)+1}^{K(l,s)+\Delta k} \text{SC}_i^l, \quad (7)$$

where  $\Delta R_l$  is computed on the current minibatch. Every  
 $T$  fine-tuning iterations, we assign the next increment  $\Delta k$   
 to the block with largest  $\Delta R_l$ , i.e., the block whose next  
 pruned tokens are most redundant.

**Timestep-wise decay.** Pruning schedule should also re-  
 flect denoising phase: Early timesteps are noise-dominated  
 and set global structure (more redundancy), while late  
 timesteps refine local details (less redundancy). Accord-  
 ingly, we adopt more aggressive pruning early and reduce  
 pruning late. In principle, one could learn a distinct pruning  
 ratio for each denoising step. However, changing  $r(l, s)$   
 across steps changes how many tokens are retained and  
 therefore modifies the execution path of MHSA and re-  
 construction. Implementing a fully step-specific schedule  
 would require instantiating (and often caching) many step-  
 dependent computation graphs and associated buffers (e.g.,  
 on mobile NPUs), which increases memory footprint and  
 engineering complexity. We therefore adopt a lightweight  
 two-phase decay that captures the major redundancy shift  
 from early (noise-dominated) to late (detail-refinement)  
 denoising, inspired by prior studies [5, 34], while requiring  
 only *two* graphs per block and keeping memory overhead

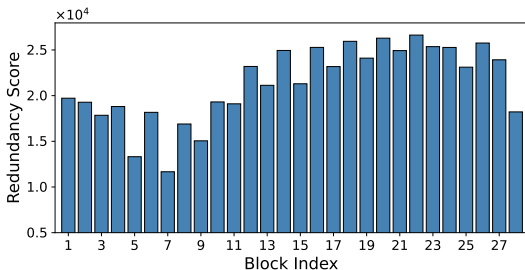
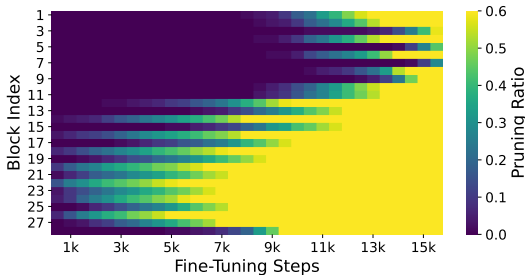
low:  $r(l, s) = \begin{cases} r(l), & 1 \leq s \leq s_0 \\ c \cdot r(l), & s_0 < s \leq S \end{cases}$  where  $r(l)$  is the  
 base pruning ratio for block  $l$ ,  $S$  is the total number of de-  
 noising steps,  $s_0$  is the phase boundary, and  $0 < c < 1$  is  
 the late-phase decay.

**Empirical study.** We demonstrate progressive pruning on  
 PixArt- $\alpha$ -1024 by increasing the pruned tokens by  $\Delta k =$   
 64 every  $T = 15$  fine-tuning steps, with a per-block cap of  
 60% tokens. Fig. 5a shows that blocks differ substantially  
 in redundancy at initialization. Consequently, the learned  
 schedule (Fig. 5b) allocates higher pruning ratios to con-  
 sistently redundant blocks while keeping low-redundancy  
 blocks lightly pruned (e.g., blocks 5 and 7). As a result,  
 CoReDiT rapidly reduces attention and end-to-end FLOPs  
 (Fig. 6a) while maintaining generation quality (Fig. 6b).

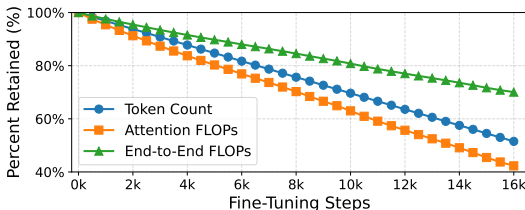
## 4. Results

### 4.1. Experimental Setup

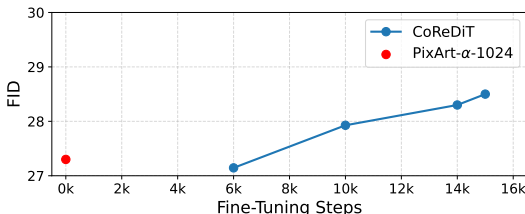
**Models and training datasets.** For text-to-image gener-  
 ation, we apply CoReDiT to the official checkpoints of  
 PixArt- $\alpha$ -1024 [6] and PixArt- $\Sigma$ -2048 [7]. We perform  
 lightweight fine-tuning with progressive pruning on the 10K  
 synthetic dataset from CLEAR [17], where images are gen-  
 erated by FLUX.1-dev [13]. For video generation, we apply  
 CoReDiT to MagicDrive-V2 [9] and fine-tune the released  
 third-stage checkpoint on nuScenes [3].

(a) Block-wise redundancy of the next  $\Delta k$  candidate pruned tokens at initialization (computed via Eq. (7)).

(b) Learned pruning ratios over fine-tuning.

Figure 5. Coherence-guided progressive pruning across transformer blocks for PixArt- $\alpha$ -1024. Higher redundancy blocks receive larger pruning ratios as fine-tuning proceeds.

(a) Efficiency during progressive pruning: retained-token ratio, attention FLOPs, and end-to-end FLOPs.



(b) Quality during progressive pruning (FID). Red dot: pretrained baseline; blue curve: CoReDiT during fine-tuning.

Figure 6. Progressive pruning dynamics for PixArt- $\alpha$ -1024: efficiency gains and quality trend over fine-tuning.

**374 Progressive pruning configuration.** For PixArt- $\alpha$ -1024,  
**375** we fine-tune with batch size of 20 on a single Nvidia H100  
**376** GPU; for PixArt- $\Sigma$ -2048, we use a total batch size of 40  
**377** across 8 Nvidia H100 GPUs. Unless otherwise stated, we

update the pruning schedule every  $T = 15$  fine-tuning  
 378 steps by pruning additional  $\Delta k = 64$  tokens in the  
 379 selected block, with a per-block pruning ratio cap of 60%;  
 380 we use a two-phase timestep schedule with decay  $c = 0.25$   
 381 and  $S - s_0 = 5$ , i.e., decay only in the final 5 denoising  
 382 steps, as detailed in Sec. 3.3. Following CLEAR [17], we  
 383 apply distillation during fine-tuning using both model out-  
 384 put loss  $L_{\text{pred}}$  and the pruned-block loss  $L_{\text{attn}}$ :  $L_{\text{distill}} =$   
 385  $L_{\text{ori}} + \alpha L_{\text{pred}} + \beta L_{\text{attn}}$ , following their hyperparameters  
 386  $\alpha = \beta = 0.5$ . For MagicDrive-V2, we prune only the  
 387 *spatial* transformer blocks, since they represent the main  
 388 computational bottleneck, with a significantly higher token  
 389 count than the temporal blocks. We fine-tune the third-stage  
 390 checkpoint with SP size of 8 across 8 Nvidia H100 GPUs.  
 391

**Evaluation metrics.** For text-to-image generation, fol-  
 392 lowing CLEAR [17], we evaluate on 10k caption-image  
 393 pairs randomly sampled from MSCOCO 2014 valida-  
 394 tion [16] using FID [10], CLIP text similarity [22], and In-  
 395 ception Score (IS) [26]. For video generation, we report  
 396 FVD [29], LPIPS [36], PSNR, and SSIM [30] for video  
 397 quality; we also evaluate conditional alignment using BEV-  
 398 Former [14] with mAP and mIoU.  
 399

## 4.2. Main Results: PixArt- $\alpha$ -1024 400

**Quality metrics.** Tab. 1 compares PixArt- $\alpha$ -1024 with  
 401 representative token-reduction and caching baselines [1, 8,  
 402 20, 37]. At  $r = 40\%$  average pruning, CoReDiT re-  
 403 duces self-attention FLOPs by 48% (24% total FLOPs)  
 404 while maintaining quality close to the baseline (FID 28.7 vs.  
 405 27.3). Notably, CLIP and IS remain competitive, indicating  
 406 that semantic alignment is preserved under substantial at-  
 407 tention compute reduction. With distillation, CoReDiT fur-  
 408 ther narrows the FID gap (27.4) with only minor changes in  
 409 IS. Fig. 7 shows that CoReDiT preserves global semantics  
 410 and fine details without introducing noticeable artifacts at  
 411 40% pruning.  
 412

**Efficiency and latency.** Tab. 2 reports both self-attention  
 413 and end-to-end latency on Nvidia H100 GPUs under two  
 414 attention backends: (i) highly-optimized kernels XFORM-  
 415 ERS.MEMORY\_EFFICIENT\_ATTENTION (Efficient Attn) and  
 416 (ii) PyTorch native attention implementations (Native Attn).  
 417 With 45% pruning, CoReDiT reduces self-attention latency  
 418 by 26% under the efficient attention backend, but end-to-  
 419 end latency improves by 11% because non-attention com-  
 420 ponents (e.g., FFN/MLP, I/O, and kernel launch overheads)  
 421 remain unchanged. Under native attention, latency im-  
 422 provements track the reduced quadratic attention cost more  
 423 closely (37% in self-attention and 25% end-to-end).  
 424

Beyond GPU inference, Fig. 8 shows per-block latency  
 425 on Qualcomm Snapdragon 8 Elite NPUs. At 50% prun-  
 426 ing, CoReDiT achieves up to  $1.72\times$  per-block speedup at  
 427

Model	FLOPs Reduction		Image Quality		
	Self-Attn	End-to-end	FID ↓	CLIP ↑	IS ↑
PixArt- $\alpha$ -1024	-	-	27.3	31.6	37.77
ToMeSD (25% ratio) [1]	-	-7%	174.6	30.2	11.68
DiffPruning [8]	-	-9%	34.6	32.0	-
EcoDiff [37]	-	-9%	32.2	32.0	-
DeepCache (N=2) [20]	-	-25%	31.6	33.1	37.44
CoReDiT ( $r = 40\%$ )	-48%	-24%	28.7	32.1	37.96
CoReDiT ( $r = 40\%$ ) + distillation	-48%	-24%	27.4	31.9	36.85
CoReDiT ( $r = 45\%$ )	-55%	-28%	29.3	31.9	36.67
CoReDiT ( $r = 45\%$ ) + distillation	-55%	-28%	28.5	31.9	36.65

Table 1. Results on PixArt- $\alpha$ -1024, where  $r$  denotes the average pruning ratio (across blocks and timesteps) in CoReDiT.

Model	FLOPs Reduction		Latency (Efficient Attn)		Latency (Native Attn)	
	Self-Attn	End-to-end	Self-Attn	End-to-end	Self-Attn	End-to-end
PixArt- $\alpha$ -1024	-	-	0.68s	1.68s	2.16s	3.17s
CoReDiT ( $r = 45\%$ )	-55%	-28%	0.50s (-26%)	1.50s (-11%)	1.36s (-37%)	2.38s (-25%)

Table 2. Comparison of FLOPs and latency on Nvidia H100 GPUs, measured with batch size 64.

1600 resolution. Notably, the baseline PixArt- $\alpha$  encounters out-of-memory (OOM) error at resolutions above 1600, whereas CoReDiT runs up to 1920 resolution and matches the baseline latency at 1600. These results indicate that structured attention skipping can provide not only runtime gains but also meaningful *memory headroom* for higher-resolution deployment.

### 4.3. High Resolution: PixArt- $\Sigma$ -2048

We apply CoReDiT to PixArt- $\Sigma$ -2048 to evaluate high-resolution image generation. As shown in Tab. 3, at an average pruning level of 23%, CoReDiT reduces self-attention FLOPs by 32% (25% total FLOPs), translating into a 22% reduction in self-attention latency and a 15% end-to-end latency reduction (batch size 32, memory-efficient attention). Quality remains high: CLIP remains the same, while FID and IS exhibit only a modest drop (26.0 vs. 28.0 and 37.49 vs. 36.52, respectively). We attribute part of this gap to a resolution-domain mismatch: our fine-tuning dataset is at 1K (upscaled to 2K) while evaluation uses native 2K images. Upscaled 1K data can under-represent 2K high-frequency details and long-range structures; fine-tuning on native 2K images is expected to narrow the remaining quality gap.

### 4.4. Video Generation: MagicDrive-V2

We extend CoReDiT, which is broadly applicable across vision tasks and agnostic to input modalities, to the state-of-the-art video generation DiT for autonomous driving - MagicDrive-V2 [9]. We prune only spatial transformer blocks and achieve a 39% self-attention FLOPs reduction (8% total FLOPs at the full-model level) as shown in Tab. 4. Despite this reduction, CoReDiT preserves strong

per-frame quality (PSNR, LPIPS, and SSIM remain close to the baseline) and maintains conditional alignment, as measured by BEVFormer mAP and mIoU.

Currently, CoReDiT is applied only to spatial transformer blocks in which each frame is processed independently. Consequently, we observed a drop in FVD, which reflects temporal consistency in the generated video. Future work will focus on integrating temporal consistency objectives, such as cross-frame coherence, to mitigate this limitation. We expect such enhancements will improve FVD performance while maintaining the computational efficiency demonstrated by our current approach.

### 4.5. Ablation Study

Tab. 5 presents our ablation studies to evaluate key components of CoReDiT on PixArt- $\alpha$ -1024 (all with distillation): (1) *Random Token Selection*. To evaluate the importance of spatial-coherence-based token selection, we conduct experiments of randomly sampling tokens to skip in each block; this yields extremely poor quality (FID 644.5 at 23% pruning), proving that spatial-coherence based selection is effective. (2) *Skipping Tokens without Reconstruction*. To illustrate the contribution of reconstruction, we disable it and simply skip selected tokens. At 30% pruning ratio, FID degrades to 30.4, indicating that reconstruction is necessary to preserve semantic fidelity. (3) *Uniform Pruning then Fine-tuning*. To justify the effectiveness of progressive pruning, we apply a uniform 50% per-block pruning to the pretrained model and then fine-tune; this results in worse FID 32.3 at average 40.6% pruning, highlighting the advantage of progressive, block-adaptive pruning.

(a) Pretrained PixArt- $\alpha$ -1024 model [6]

(b) Proposed CoReDiT with 40% FLOPs Pruning. The semantics are preserved and no additional artifacts.

Figure 7. Qualitative comparison to PixArt- $\alpha$ -1024.

Model (ratio)	FLOPs Reduction		Latency		Image Quality		
	Self-Attn	End-to-end	Self-Attn	End-to-end	FID ↓	CLIP ↑	IS ↑
PixArt- $\Sigma$ -2048	-	-	4.62s	6.68s	26.0	31.4	37.49
CoReDiT ( $r = 23\%$ ) + distillation	-32%	-25%	3.61s (-22%)	5.67s (-15%)	28.0	31.4	36.52

Table 3. Results on PixArt- $\Sigma$ -2048, with latency measured on Nvidia H100 GPUs using Efficient Attention.

Model	FLOPs Reduction		Video Quality				Cond. Alignment	
	Self-Attn	End-to-end	PSNR ↑	LPIPS ↓	SSIM ↑	FVD ↓	mAP ↑	mIoU ↑
MagicDrive-V2	-	-	14.28	0.422	0.372	107.8	18.4%	21.5%
CoReDiT ( $r = 26\%$ )	-39%	-8%	14.25	0.424	0.378	119.8	18.1%	21.2%

Table 4. Video generation results on video quality and conditional alignment via 3D object detection

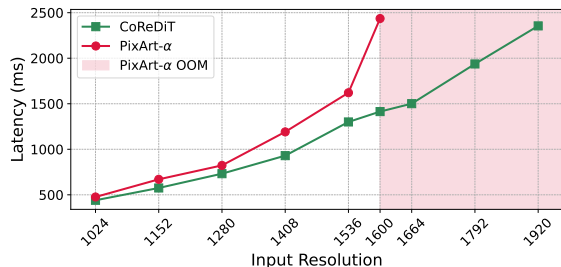


Figure 8. Per-block latency comparison on Qualcomm Snapdragon 8 Elite NPUs. OOM: out-of-memory.

Experiment	Image Quality		
	FID ↓	CLIP ↑	IS ↑
PixArt- $\alpha$ -1024	27.3	31.6	37.77
CoReDiT ( $r = 45\%$ )	28.5	31.9	36.65
Random selection ( $r = 23\%$ )	644.5	21.5	1.00
Disable reconstruction ( $r = 30\%$ )	30.4	31.7	34.20
Uniform ratio ( $r = 41\%$ )	32.3	31.5	33.68

Table 5. Ablation study on PixArt- $\alpha$ -1024, with distillation in all experiments.

489

## 5. Conclusion

490

491

492

493

494

495

496

497

We introduce CoReDiT, a token pruning framework for DiTs that exploits spatial coherence as a practical signal of local redundancy. CoReDiT combines (1) an efficient pre-attention selector that skips high-coherence tokens, (2) a coherence-guided reconstruction operator that preserves a dense token lattice and mitigates pruning artifacts, and (3) a progressive, block-adaptive pruning schedule that allocates pruning where redundancy is highest across blocks

and denoising steps. Experiments on state-of-the-art image and video generators demonstrate substantial reductions in self-attention compute and measurable end-to-end speedups on both cloud GPUs and mobile NPUs, while maintaining strong perceptual quality and conditional alignment. In future work, we plan to (1) further reduce non-attention overhead by extending efficiency to FFN/MLP components, and (2) improve video temporal consistency by incorporating cross-frame coherence objectives and pruning strategies that explicitly model temporal redundancy.

498

499

500

501

502

503

504

505

506

507

508

## References

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

- [1] 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. 2, 6, 7
- [2] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token merging: Your vit but faster. *arXiv preprint arXiv:2210.09461*, 2022. 1, 3
- [3] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multi-modal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020. 5
- [4] Keshigeyan Chandrasegaran, Michael Poli, Daniel Y Fu, Dongjun Kim, Lea M Hadzic, Manling Li, Agrim Gupta, Stefano Massaroli, Azalia Mirhoseini, Juan Carlos Niebles, et al. Exploring diffusion transformer designs via grafting. *arXiv preprint arXiv:2506.05340*, 2025. 2, 4
- [5] Shuning Chang, Pichao Wang, Jiasheng Tang, and Yi Yang. Flexdit: Dynamic token density control for diffusion transformer. *arXiv preprint arXiv:2412.06028*, 2024. 2, 4, 5
- [6] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart- $\alpha$ : Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023. 5, 8
- [7] 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*, pages 74–91. Springer, 2024. 5
- [8] Gongfan Fang, Xinyin Ma, and Xinchao Wang. Structural pruning for diffusion models. In *Advances in Neural Information Processing Systems*, 2023. 2, 6, 7
- [9] Ruiyuan Gao, Kai Chen, Bo Xiao, Lanqing Hong, Zhenguo Li, and Qiang Xu. Magicdrive-v2: High-resolution long video generation for autonomous driving with adaptive control. *arXiv preprint arXiv:2411.13807*, 2024. 5, 7
- [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. *Advances in neural information processing systems*, 30, 2017. 6
- [11] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2
- [12] Sehoon Kim, Sheng Shen, David Thorsley, Amir Gholami, Woosuk Kwon, Joseph Hassoun, and Kurt Keutzer. Learned token pruning for transformers. In *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*, pages 784–794, 2022. 1, 3
- [13] Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024. 5
- [14] Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai. Bevformer: Learning bird’s-eye-view representation from multi-camera images via spatiotemporal transformers. *arXiv preprint arXiv:2203.17270*, 2022. 6
- [15] Youwei Liang, Chongjian Ge, Zhan Tong, Yibing Song, Jue Wang, and Pengtao Xie. Not all patches are what you need: Expediting vision transformers via token reorganizations. *arXiv preprint arXiv:2202.07800*, 2022. 3
- [16] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 6
- [17] Songhua Liu, Zhenxiong Tan, and Xinchao Wang. Clear: Conv-like linearization revs pre-trained diffusion transformers up. *arXiv preprint arXiv:2412.16112*, 2024. 5, 6
- [18] Wenbo Lu, Shaoyi Zheng, Yuxuan Xia, and Shengjie Wang. Toma: Token merge with attention for diffusion models. *arXiv preprint arXiv:2509.10918*, 2025. 2
- [19] Zhengyao Lv, Chenyang Si, Junhao Song, Zhenyu Yang, Yu Qiao, Ziwei Liu, and Kwan-Yee K Wong. Fastercache: Training-free video diffusion model acceleration with high quality. *arXiv preprint arXiv:2410.19355*, 2024. 2
- [20] 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. 6, 7
- [21] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023. 1, 2
- [22] 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 *International conference on machine learning*, pages 8748–8763. PmLR, 2021. 6
- [23] Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dynamicvit: Efficient vision transformers with dynamic token sparsification. *Advances in neural information processing systems*, 34:13937–13949, 2021. 1
- [24] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2
- [25] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 2
- [26] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, et al. Improved techniques for training gans. *Advances in neural information processing systems*, 29, 2016. 6

- 620 [27] Manish Kumar Singh, Rajeev Yasarla, Hong Cai, Mingu  
621 Lee, and Fatih Porikli. Tosa: Token selective atten-  
622 tion for efficient vision transformers. *arXiv preprint*  
623 *arXiv:2406.08816*, 2024. 1, 3
- 624 [28] Ethan Smith, Nayan Saxena, and Aninda Saha. Todo: Token  
625 downsampling for efficient generation of high-resolution im-  
626 ages. *arXiv preprint arXiv:2402.13573*, 2024. 2
- 627 [29] Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach,  
628 Raphael Marinier, Marcin Michalski, and Sylvain Gelly. To-  
629 wards accurate generative models of video: A new metric &  
630 challenges. *arXiv preprint arXiv:1812.01717*, 2018. 6
- 631 [30] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Si-  
632 moncelli. Image quality assessment: from error visibility to  
633 structural similarity. *IEEE transactions on image processing*,  
634 13(4):600–612, 2004. 6
- 635 [31] Felix Wimbauer, Bichen Wu, Edgar Schoenfeld, et al. Cache  
636 me if you can: Accelerating diffusion models through block  
637 caching. In *Proceedings of the IEEE/CVF Conference on*  
638 *Computer Vision and Pattern Recognition*, pages 6211–  
639 6220, 2024. 2
- 640 [32] Haocheng Xi, Shuo Yang, Yilong Zhao, Chenfeng Xu,  
641 Muyang Li, Xiuyu Li, Yujun Lin, Han Cai, Jintao Zhang,  
642 Dacheng Li, et al. Sparse videogen: Accelerating video  
643 diffusion transformers with spatial-temporal sparsity. *arXiv*  
644 *preprint arXiv:2502.01776*, 2025. 2
- 645 [33] Huanrui Yang, Hongxu Yin, Maying Shen, Pavlo  
646 Molchanov, Hai Li, and Jan Kautz. Global vision trans-  
647 former pruning with hessian-aware saliency. In *Proceedings*  
648 *of the IEEE/CVF conference on computer vision and pattern*  
649 *recognition*, pages 18547–18557, 2023. 5
- 650 [34] Haoran You, Connelly Barnes, Yuqian Zhou, Yan Kang,  
651 Zhenbang Du, Wei Zhou, Lingzhi Zhang, Yotam Nitzan, Xi-  
652 aoyang Liu, Zhe Lin, et al. Layer-and timestep-adaptive dif-  
653 ferentiable token compression ratios for efficient diffusion  
654 transformers. In *Proceedings of the Computer Vision and*  
655 *Pattern Recognition Conference*, pages 18072–18082, 2025.  
656 1, 2, 4, 5
- 657 [35] Peiyuan Zhang, Yongqi Chen, Haofeng Huang, Will Lin,  
658 Zhengzhong Liu, Ion Stoica, Eric Xing, and Hao Zhang.  
659 Vsa: Faster video diffusion with trainable sparse attention.  
660 *arXiv preprint arXiv:2505.13389*, 2025. 2
- 661 [36] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shecht-  
662 man, and Oliver Wang. The unreasonable effectiveness of  
663 deep features as a perceptual metric. In *Proceedings of the*  
664 *IEEE conference on computer vision and pattern recogni-*  
665 *tion*, pages 586–595, 2018. 6
- 666 [37] Yang Zhang, Er Jin, Yanfei Dong, Ashkan Khakzar, Philip  
667 Torr, Johannes Stegmaier, and Kenji Kawaguchi. Effort-  
668 less efficiency: Low-cost pruning of diffusion models. *arXiv*  
669 *preprint arXiv:2412.02852*, 2024. 6, 7
- 670 [38] Wangbo Zhao, Yizeng Han, Jiasheng Tang, Kai Wang, Yib-  
671 ing Song, Gao Huang, Fan Wang, and Yang You. Dynamic  
672 diffusion transformer. *ICLR*, 2025. 2
- 673 [39] Chang Zou, Xuyang Liu, Ting Liu, Siteng Huang, and Lin-  
674 feng Zhang. Accelerating diffusion transformers with token-  
675 wise feature caching. *arXiv preprint arXiv:2410.05317*,  
676 2024. 2