# DraftAttention: Fast Video Diffusion via Low-Resolution Attention Guidance

**Xuan Shen**[1*], **Chenxia Han**[2*], **Yufa Zhou**[3*], **Yanyue Xie**[1], **Yifan Gong**[4],
**Quanyi Wang**[5], **Yiwei Wang**[6], **Yanzhi Wang**[1†], **Pu Zhao**[1†], **Jiuxiang Gu**[4†]

[1]Northeastern University, [2]CUHK, [3]Duke University,
[4]Adobe Research, [5]NUIST, [6]UCM
shen.xu@northeastern.edu, cxhan@cse.cuhk.edu.hk

## Abstract

Diffusion transformer–based video generation models (DiTs) have recently attracted widespread attention for their excellent generation quality. However, their computational cost remains a major bottleneck—attention alone accounts for over 80% of total latency, and generating just 8 seconds of 720p video takes tens of minutes—posing serious challenges to practical application and scalability. To address this, we propose the `DraftAttention`, a training-free framework for the acceleration of video diffusion transformers with dynamic sparse attention on GPUs. We apply down-sampling to each feature map across frames in the compressed latent space, enabling a higher-level receptive field over the latent composed of hundreds of thousands of tokens. The low-resolution draft attention map, derived from draft query and key, exposes redundancy both spatially within each feature map and temporally across frames. We reorder the query, key, and value based on the draft attention map to guide the sparse attention computation in full resolution, and subsequently restore their original order after the attention computation. This reordering enables structured sparsity that aligns with hardware-optimized execution. Our theoretical analysis demonstrates that the low-resolution draft attention closely approximates the full attention, providing reliable guidance for constructing accurate sparse attention. Experimental results show that our method outperforms existing sparse attention approaches in video generation quality and achieves up to $1.75\times$ end-to-end speedup on GPUs. Code: https://github.com/shawnricecake/draft-attention

## 1 Introduction

Diffusion Transformers (DiTs) [1] have emerged as a powerful paradigm for visual generative tasks across both image and video generation, surpassing the traditional UNets [2].

Video generation with DiTs adopts spatiotemporal 3D full attention to extend image-based generation to the temporal domain [3], leading to visually coherent high-quality video generation performance [4, 5, 6], validating the effectiveness of DiTs for video generation. Despite the superior generation performance with DiTs, it remains computationally expensive due to the attention mechanism in transformers. The quadratic complexity with respect to context length [7] becomes a significant computational
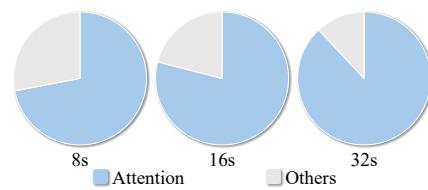


Figure 1: FLOPs breakdown for 720p video generation with Hunyuan Video.

---

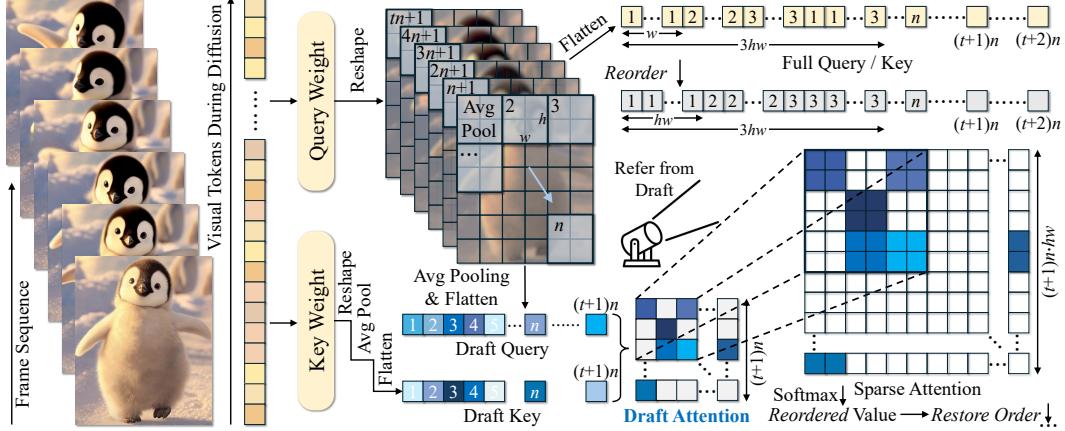[*]Equal Contribution
[†]Corresponding Authors

Figure 2: Whole `DraftAttention` Pipeline. Both the query and key are reshaped into sequences of feature maps across frames, then downsampled via average pooling to produce the low-resolution draft query and draft key. Draft attention is computed using the flattened draft query and key. The full-resolution query and key need to be reordered for the alignment of draft attention guidance.

bottleneck when handling sequences with hundreds of thousands of tokens. For example, as shown in Figure 1, the Hunyuan Video model [5] spends over 80% of its total computation on the attention mechanism when generating videos longer than 16 seconds. As a result, the slow generation speed limits the application and deployment of these promising video generation models across a range of practical tasks.

Fortunately, pioneering works [8, 9, 10, 11] on Large Language Models (LLMs) [12, 13, 14, 15] has demonstrated substantial redundancy in the attention mechanism, offering an opportunity for acceleration by introducing sparsity into the attention. Inspired by this, recent works [16, 17] explore the sparse attention methods for video generation models, demonstrating promising speedups while preserving generation quality. Specifically, two static sparse attention patterns (targeting spatial and temporal dimensions respectively) are explored in Sparse VideoGen [16] to reduce redundancy, with relatively significant performance degradation under large sparsity because of non-adaptive static patterns. To mitigate this issue, dynamic sparse attention is investigated in AdaSpa [17] to perform full attention once for different prompts as a warm-up to guide subsequent sparsity. Although AdaSpa provides prompt-dependent sparse patterns, patterns still remains static during the diffusion process.

**Framework Overview.** Motivated by the absence of true dynamic sparse attention at the per-module level, we investigate a more fine-grained design—adapting the sparse attention patterns dynamically for each specific attention module. In this paper, we propose an efficient sparse attention method, `DraftAttention`, as shown in Figure 2, which leverages draft attention to dynamically generate a sparse pattern for each attention module, enabling efficient acceleration of video diffusion transformers. The key idea is to compute the draft attention based on downsampled low-resolution query and key, thus identifying the most important areas in the attention map with minor computational overhead. The resulting low-resolution sparse mask then guides full-resolution sparse attention, with effective reordering applied to ensure fast, hardware-friendly execution.

**Great Advantages.** We highlight the following advantages with our draft attention method: **(i)** (Efficiency) The computation of draft attention map is lightweight, as it operates on a reduced number of tokens, thereby lowering the quadratic complexity of the attention mechanism. **(ii)** (Effectiveness) The draft attention captures high-level representations and preserves essential visual patterns for videos, leading to an effective mask to identify the critical structures in attention mechanism. **(iii)** (Plug-and-Play) Our method requires no additional training and integrates seamlessly as a plug-and-play module into existing video diffusion transformers for handling long input sequences.

**Theoretical Justification.** We also present the theoretical analysis that formally characterizes how the low-resolution draft attention effectively guides the full-resolution attention mechanism. Specifically, we show that the upper bound of the difference between the full-resolution attention map and the draft attention map remains controlled. Meanwhile, we show that the error introduced by the sparse pattern derived from the draft attention map remains bounded.

**Hardware Friendliness.** To align the region-level sparsity with token-level computations, we apply a deterministic reordering of tokens such that entries in each region become contiguous in memory, ensuring hardware-friendly execution of sparse attention.

**Comprehensive Experiments.** In our experiments, we use an $8\times16$ pooling kernel with a stride equal to the kernel size, reducing the number of tokens by a factor of 128. This configuration also matches the efficient block size supported by efficient attention computation frameworks [7, 18]. Meanwhile, through reordering, we group the scattered sparse patterns into a contiguous format, allowing 128 visual tokens within each kernel to be processed in a single stage—either computed or skipped. This enables both accurate and faster sparse attention at full resolution. Such aggressive downsampling also incurs minimal computational overhead for the low-resolution draft attention. Meanwhile, our method outperforms other sparse attention methods on video generation tasks across various resolutions under the same computational budget. It achieves up to a $1.75\times$ end-to-end speedup on GPUs, demonstrating strong practical efficiency and scalability for long video sequences without compromising generation quality. Our contributions are summarized as follows,

**1.** We introduce a vision-centric perspective on spatial and temporal redundancy in video diffusion, using pooling to extract high-level representations with a broader receptive field. Building on this, we propose `DraftAttention`, a hardware-friendly approach that accelerates video diffusion transformers using guidance from low-resolution draft attention.

**2.** We provide a theoretical analysis demonstrating the controlled difference between full-resolution attention and low-resolution draft attention, as well as the bounded error introduced by the sparse pattern derived from the draft attention map, thereby justifying the effectiveness of our design.

**3.** Experimental results show that `DraftAttention` achieves better video generation quality compared to other sparse attention methods with same computation cost. Meanwhile, on GPUs, our method achieves up to $1.75\times$ end-to-end acceleration for video generation.

## 2 Related Works

### 2.1 Efficient Diffusion Models

**Diffusion Model Compression.** Weight quantization is a common approach to compress diffusion models and achieve acceleration [19]. Previous works [20, 21, 22] propose optimal quantization methods to quantize attention weights to INT8, INT4/FP8, or even FP4, which achieve high compression ratios for the diffusion model size. Also, other works explore efficient architectures [23] including linear attention or high-compression auto-encoders [24] to accelerate the diffusion and improve model performance, which extends the scalability of diffusion models. Our method is orthogonal to these techniques and integrates with them to yield additional performance gains.

**Reduce Diffusion Steps.** Some distillation-based works [25, 26] adopt training for the simpler to build few-step diffusion models, which accelerates the diffusion progress by reducing the steps. However, such distillation techniques require expensive re-training or fine-tuning, which is impractical for the application of most video diffusion models. In contrast, our approach directly uses off-the-shelf pre-trained models without any additional training.

### 2.2 Sparse Attention Methods

Attention mechanisms exhibit inherent sparsity [27], allowing computational acceleration by limiting interactions to a subset of the key-value pair. StreamingLLM [10] explores the temporal locality with attention sinks to further preserve sparse attention model performance. H2O [8] identifies a small set of Heavy Hitter tokens that dominate overall attention scores. DuoAttention [28] and MInference [11] demonstrate distinct sparse patterns across different attention heads. XAttention [29] leverages the sum of antidiagonal values in the attention matrix to provide a powerful proxy for block importance, resulting in high sparsity and dramatically accelerated inference. Sparse VideoGen [16] explores spatial and temporal heads in video diffusion models to improve the inference efficiency. AdaSpa [17] applies dynamic block-sparse masking with online token importance search, accelerating video diffusion without fine-tuning. These works collectively show that such transformer-based models contain significant redundancy in their attention mechanisms. This motivates our exploration of dynamic, fine-grained sparse attention patterns for video diffusion transformers.

# 3 Methodology

We introduce the framework of our draft attention in great detail to first identify critical areas in draft attention with a low-resolution mask and then apply the mask to full-resolution attention. Next theoretical analysis for the draft attention and the corresponding sparse attention is presented to demonstrate the effectiveness of our design. Moreover, we provide a deterministic reordering of tokens to align the region-level sparsity with token-level computation, ensuring efficient hardware-friendly execution.

## 3.1 Draft Attention

Full attention over long video sequences is prohibitively expensive due to its quadratic complexity in sequence length. However, many interactions in video are spatially and temporally localized. We leverage this structure by introducing a two-stage attention mechanism: a lightweight *draft attention* phase that estimates regional relevance, followed by a masked sparse attention applied to the full-resolution sequence.

We first define the full attention computation below.

**Definition 3.1** (Full Attention). *Given hidden states $X \in \mathbb{R}^{n \times d}$, the full attention output is:*

$$\mathsf{Attn}(X) = \mathsf{Softmax}\left(\frac{QK^\top}{\sqrt{d}}\right) V \in \mathbb{R}^{n \times d}, \tag{1}$$

*where $Q = XW_Q$, $K = XW_K$, $V = XW_V$ are the query, key, and value projections, and $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$ are learned weight matrices.*

To reduce computation, we downsample $Q$ and $K$ via average pooling, forming a low-resolution draft attention map to guide sparsity.

**Definition 3.2** (Draft Attention via Average Pooling). *Given hidden states $X \in \mathbb{R}^{n \times d}$, representing spatial-temporal tokens across frames, we partition the sequence into $g \ll n$ disjoint regions $\{R_i\}_{i=1}^g$, where each region $R_i \subset [n]$ corresponds to a pooled spatial patch over time. Each $R_i$ is an unordered set of token indices. Let $Q$ and $K$ be the projected queries and keys. The draft query and draft key representations are obtained by average pooling over each region:*

$$\widetilde{Q}_i = \frac{1}{|R_i|} \sum_{j \in R_i} Q_j, \quad \widetilde{K}_i = \frac{1}{|R_i|} \sum_{j \in R_i} K_j, \quad for \ i = 1, \ldots, g. \tag{2}$$

*The resulting low-resolution **draft attention map** is computed as:*

$$A_{\mathrm{draft}} = \mathsf{Softmax}\left(\frac{\widetilde{Q}\widetilde{K}^\top}{\sqrt{d}}\right) \in \mathbb{R}^{g \times g}. \tag{3}$$

*This map approximates region-level relevance and is used to guide sparse attention over the full-resolution sequence.*

The computation cost of the low-resolution draft attention map is minor compared with the full-resolution attention computation, as it operates on a reduced number of tokens and thereby lowers the quadratic complexity of the attention mechanism.

**Guided Sparsity via Draft Attention.** To reduce the cost of full attention, we extract a structured sparsity pattern from the draft attention map $A_{\mathrm{draft}} \in \mathbb{R}^{g \times g}$ by retaining only a fraction $r \in (0, 1)$ of the most salient region-to-region interactions. We define a binary mask $M \in \{0, 1\}^{g \times g}$, where $M_{ij} = 1$ indicates that region $R_i$ is permitted to attend to region $R_j$, and $M_{ij} = 0$ otherwise. The mask is constructed by selecting the top-scoring entries in $A_{\mathrm{draft}}$ under a fixed sparsity ratio $r$.

To align the region-level sparsity with token-level computation, we apply a deterministic reordering of tokens such that entries in each region $R_i$ become contiguous. This facilitates efficient masking and block-wise computation in sparse attention. We provide more details for reordering in Section 3.3.

4

This region-level sparsity pattern is then lifted to token resolution by defining a full-resolution binary mask $\widehat{M} \in \{0, 1\}^{n \times n}$:

$$\widehat{M}_{uv} = M_{ij} \quad \text{if } u \in R_i, \ v \in R_j. \tag{4}$$

In general, the attention map is split into multiple non-overlapping regions by the pooling kernels. For each region, all its elements are either computed for attention or skipped for acceleration. The determination for whether to skip each region is denoted by the low-resolution binary mask $M$ for all regions, with $\widehat{M}$ as its full-resolution mask for all elements (i.e., tokens).

Sparse attention is then computed by applying the mask to the full attention scores:

$$\mathsf{SparseAttn}(X) = \mathsf{Softmax} \left( \left( \frac{QK^\top}{\sqrt{d}} \right) \odot \widehat{M} \right) V, \tag{5}$$

where $\odot$ denotes element-wise/Hadamard product. This formulation retains the most relevant interactions while enforcing structured sparsity for improved computational efficiency.

## 3.2 Theoretical Analysis

We present Frobenius-norm bounds quantifying the error introduced by our two-stage approximation strategy: (1) average pooling (draft attention), and (2) structured sparsification via top-$r$ indexing.

### 3.2.1 Error from Draft Attention

Let the input sequence be partitioned into $g$ disjoint regions $\{R_i\}_{i=1}^{g}$ of equal size $|R_i| = n/g$. Define the full-resolution attention logits and their pooled approximation as:

$$S_{uv} := \langle Q_u, K_v \rangle, \quad \widetilde{S}_{ij} := \langle \widetilde{Q}_i, \widetilde{K}_j \rangle, \qquad u, v \in [n], \ i, j \in [g], \tag{6}$$

where $\widetilde{Q}_i = \frac{1}{|R_i|} \sum_{u \in R_i} Q_u$ and similarly for $\widetilde{K}_j$.

We restore the region-level scores $\widetilde{S} \in \mathbb{R}^{g \times g}$ to full resolution by defining a block-constant approximation:

$$(S_{\mathrm{draft}})_{uv} := \widetilde{S}_{ij} \quad \text{for } u \in R_i, \ v \in R_j. \tag{7}$$

Define the worst-case deviation between token-level logits and their region-averaged counterpart as:

$$\delta := \max_{i,j} \ \max_{u \in R_i, v \in R_j} \left| S_{uv} - \widetilde{S}_{ij} \right|. \tag{8}$$

**Theorem 3.3** (Draft Attention Error). *If all regions have equal size $|R_i| = n/g$, then the Frobenius-norm error between the full and draft logit matrices is bounded by:*

$$\| S - S_{\mathrm{draft}} \|_F \leq \delta\, n. \tag{9}$$

The detailed proof is shown in Appendix A.

**Remark 3.4.** *Theorem 3.3 quantifies the approximation error introduced by replacing token-level attention logits with block-wise averages obtained via average pooling. In practice, if tokens within a region are similar—such as in videos with local temporal consistency or spatial smoothness—the difference $|S_{uv} - \widetilde{S}_{ij}|$ remains small for most $(u, v)$. Consequently, the overall Frobenius-norm error $\| S - S_{\mathrm{draft}} \|_F$ scales with a modest $\delta$, leading to minimal distortion in the attention structure. This justifies using the low-resolution draft map as a proxy for full-resolution attention in computationally constrained settings.*

### 3.2.2 Error from Sparsity Mask

We now consider the additional error introduced by sparsifying the logits based on the top-$r$ draft attention values. Let $\widetilde{S}_{(1)} \geq \cdots \geq \widetilde{S}_{(g^2)}$ be the sorted region-level scores. Define the threshold:

$$t := \widetilde{S}_{(\lceil rg^2 \rceil)}, \tag{10}$$

and let $M_{ij} = 1$ if $\widetilde{S}_{ij} \geq t$ and 0 otherwise. The mask is lifted to token resolution by $\widehat{M}_{uv} = M_{ij}$ for $u \in R_i, \ v \in R_j$.
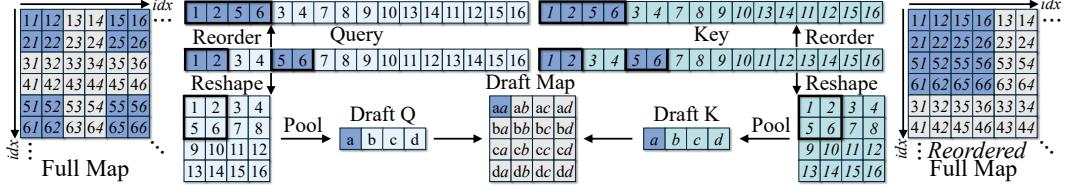
Figure 3: Illustration for the necessity of the reordering. The "xy" in attention map denotes attentivity between token x in query and token $y$ in key. Grouping the sparse pattern enables hardware-friendly layout, leading to faster attention computation.

**Theorem 3.5** (Sparsity Mask Error). *Under uniform region size $|R_i| = n/g$, the error from masking the logits satisfies:*

$$\|S - S \odot \widehat{M}\|_F \le n(\delta + t)\sqrt{1 - r}. \tag{11}$$

The detailed proof is shown in Appendix A.

**Remark 3.6.** *Theorem 3.5 captures the additional error introduced by enforcing structured sparsity through top-$r$ indexing. This error is controlled by the threshold value $t$, which defines the weakest retained region-level interaction. When the draft attention distribution is concentrated—i.e., a small fraction of regions account for most of the mass—top-$r$ masking retains the most informative blocks while discarding only low-scoring interactions, keeping $(\delta + t)$ small. This ensures that sparsification preserves the dominant attention patterns with bounded deviation.*

Together, Theorems 3.3 and 3.5 provide a principled decomposition of the total approximation error: one from average pooling, and one from sparsity. Their combined bound shows that draft attention is an efficient surrogate for full attention, maintaining structural fidelity while enabling substantial computational savings. This justifies its use in long-context video diffusion transformers, where local smoothness and sparse relevance patterns are common.

### 3.3 Reordering for Patch-Aligned Sparse Attention

To enable accurate and efficient sparse attention that respects spatial structure, we apply a deterministic reordering algorithm (Algorithm 1) to the flattened full-resolution token sequence. As shown in Figure 3, the goal is to align the memory layout of full-resolution tokens with the spatial region structure used in low-resolution draft attention. This alignment ensures that the region-level sparsity patterns are directly and efficiently propagated to full-resolution attention through block-level masking.

---

**Algorithm 1:** Generate Reorder Index

**Input:** Frame size $(H, W)$, patch size $(h, w)$, number of frames $F$
**Output:** Permutation $\pi \in [n]$ where $n = F \cdot H \cdot W$

$\pi \leftarrow [\,]$;
**for** $f = 0$ **to** $F - 1$ **do**
   **for** $i = 0$ **to** $H/h - 1$ **do**
      **for** $j = 0$ **to** $W/w - 1$ **do**
         **for** $u = 0$ **to** $h - 1$ **do**
            **for** $v = 0$ **to** $w - 1$ **do**
               $y \leftarrow i \cdot h + u, x \leftarrow j \cdot w + v$;
               idx $\leftarrow f \cdot H \cdot W + y \cdot W + x$;
               Append idx to $\pi$;

**return** $\pi$

---

**Justification.** In the default row-major layout, spatial tokens are appended row-wise within each frame, causing spatial patches to be scattered in memory. This fragmentation hinders efficient usage of sparse attention kernels, which rely on contiguous blocks in fixed size for the optimal performance. As illustrated in Figure 3, tokens 1, 2, 5, and 6 are spatial neighbors but are not stored consecutively in the memory of full attention map (i.e., left side of Figure 3) due to the presence of tokens 3 and 4. While it is still possible to gather these tokens and compute their average, this process is highly inefficient. Similarly, masking out these scattered blocks is also inefficient, as it effectively reduces the block size, which in turn lowers arithmetic intensity, causes uncoalesced memory access, and increases the number of kernel launches.

**Design.** We divide each frame into non-overlapping patches of size $h \times w$. For each frame, tokens within the same patch are grouped contiguously. Unlike prior methods (e.g., SVG [16]) that overlook misalignment issues when the kernel size does not divide evenly into the latent feature map size, our per-frame design preserves the completeness of each feature map, generating more reliable captured high-level representations. Meanwhile, this per-frame design ensures that each patch in a frame is stored as a contiguous block, matching the structure of the downsampled low-resolution queries and keys used in draft attention. For instance, tokens 1, 2, 5, and 6 belong to the same patch and are reordered to appear consecutively in both the query and key sequences, as illustrated at the top of Figure 3. This reordering ensures that each entry in draft attention map (e.g., a$a$) corresponds to a specific block ($\{1, 2, 5, 6\}$ from query and $\{1, 2, 5, 6\}$ from key) within reordered full attention map.

**Execution.** Applying the permutation $\pi$ ensures that tokens grouped in each $h \times w$ patch are stored contiguously in memory, enabling efficient block-wise indexing and masking. This structured layout aligns the memory access pattern with the computational needs of sparse attention operations. This is especially critical for efficient execution with frameworks like FlashAttention [7] and Block Sparse Attention [18], which leverage fused GPU kernels that operate on fixed-size blocks.

**Restoration.** After sparse attention is applied in the reordered space (i.e., the attention computation for reordered query, key, and value), we apply the inverse permutation $\pi^{-1}$ (Algorithm 2) to restore the original spatial-temporal layout for the following correct model inference.

**Benefit.** This reordering bridges the gap between the coarse-grained sparsity structure derived from draft attention and the fine-grained full-resolution attention computation. This layout guarantees that pooled regions align cleanly with memory blocks, preserving spatial locality and enabling predictable, coalesced memory access. As a result, it supports efficient masking and ensures compatibility with high-throughput attention kernels. This design significantly reduces overhead and maximizes hardware efficiency during attention computation.

---

**Algorithm 2:** Generate Restore Index

---
**Input:** Permutation $\pi \in [n]$
**Output:** Inverse permutation $\pi^{-1}$
Initialize $\pi^{-1} \leftarrow$ zero array of length $n$;
**for** $i = 0$ **to** $n - 1$ **do**
$\quad \lfloor \ \pi^{-1}_{\pi_i} \leftarrow i;$
**return** $\pi^{-1}$

---

# 4 Experimental Results

## 4.1 Experiment Setup

**Model Family.** We adopt open-sourced state-of-the-art video generation models in our experiments, including HunyuanVideo-T2V [5] for 768p resolution with 128 frames and Wan2.1-T2V [6] for both 512p and 768p resolutions with 80 frames. We use 512p and 768p resolutions to align with the $8 \times 16$ average pooling kernel (with stride equal to the kernel size), enabling convenient and consistent downsampling of visual tokens during the diffusion process. This is because the corresponding latent sizes—$32 \times 48$ for 512p and $48 \times 80$ for 768p—are perfectly divisible by the $8 \times 16$ kernel, ensuring efficient and artifact-free pooling. Note that our method supports video generation at any resolution by applying appropriate padding. Following prior works [16, 30, 31, 32], we retain full attention across all methods for the first 25% of denoising steps to preserve the video generation quality. We adopt Block Sparse Attention [18] for the implementation of our method and mainly compare our method with the Sparse VideoGen (SVG) [16]. We observe discrepancies in the generation results of the Wan2.1-T2V model between our method and SVG, due to difference of codebases. To ensure a fair comparison, we provide results using full attention for both methods.

**Metrics and Prompts.** We evaluate the quality of generated videos with VBench [33], and the similarity of generated videos with metrics including Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index Measure (SSIM), and Learned Perceptual Image Patch Similarity (LPIPS) [34]. Especially, we report the image quality, subject consistency, background consistency, dynamic degree, and aesthetic quality from VBench for our generated videos. All videos are generated with the prompts from the Penguin Video Benchmark [5] released by HunyuanVideo. The reported computation cost in PFLOPs includes the main diffusion transformer models, and the latency results are all tested on the H100 GPU.

Table 1: Main results of the proposed method compared to the Sparse VideoGen (SVG) [16].

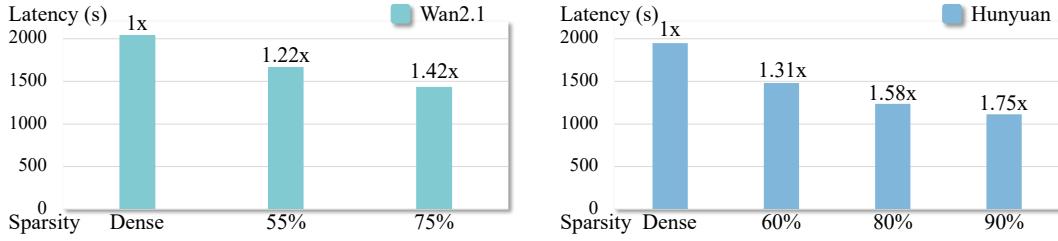| Model | Method | Sparse Ratio | PSNR ↑ | SSIM ↑ | LPIPS ↓ | Img. Qual. | Sub. Cons. | Bakg. Cons. | Dyn. Deg. | Aes. Qual. | PFLOPs ↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Wan2.1 (512p) | SVG | 0% | / | / | / | 65.1% | 95.0% | 95.9% | 44.7% | 58.9% | 145.65 |
| | | 55% | 25.61 | 83.63 | 10.42 | 65.2% | 94.8% | 95.9% | 45.2% | 58.9% | 99.26 |
| | | 75% | 23.66 | 78.80 | 15.05 | 64.7% | 94.5% | 95.7% | 45.7% | 58.6% | 91.12 |
| | Ours | 0% | / | / | / | 69.3% | 95.5% | 96.7% | 47.6% | 61.5% | 145.65 |
| | | 55% | 25.13 | **84.77** | **8.43** | 69.2% | 95.5% | 96.6% | 47.6% | 61.5% | 99.26 |
| | | 75% | 23.10 | **79.07** | **12.37** | 69.0% | 95.4% | 96.5% | 46.9% | 61.5% | 91.12 |
| Wan2.1 (768p) | SVG | 0% | / | / | / | 67.7% | 95.3% | 96.4% | 43.4% | 60.4% | 609.52 |
| | | 55% | 26.01 | 84.81 | 10.89 | 67.9% | 95.1% | 96.3% | 42.1% | 60.0% | 354.68 |
| | | 75% | 23.62 | 79.05 | 17.57 | 67.5% | 94.8% | 96.1% | 42.1% | 58.8% | 309.95 |
| | Ours | 0% | / | / | / | 67.5% | 95.7% | 97.1% | 37.7% | 60.8% | 609.52 |
| | | 55% | **29.22** | **92.16** | **5.82** | 67.4% | 95.6% | 97.0% | 37.2% | 60.8% | 354.69 |
| | | 75% | **27.17** | **88.97** | **8.71** | 67.2% | 95.6% | 97.0% | 38.6% | 60.7% | 309.95 |
| Hunyuan (768p) | Dense | 0% | / | / | / | 66.4% | 96.0% | 97.0% | 36.4% | 58.6% | 682.67 |
| | SVG | 60% | 25.80 | 84.46 | 14.20 | 66.4% | 95.9% | 97.0% | 36.6% | 58.2% | 343.72 |
| | | 80% | 24.70 | 81.90 | 17.55 | 66.0% | 95.7% | 96.9% | 33.9% | 58.1% | 295.30 |
| | | 90% | 23.48 | 78.57 | 22.60 | 65.1% | 95.4% | 96.7% | 32.8% | 57.5% | 283.20 |
| | Ours | 60% | **32.08** | **93.21** | **5.58** | **66.4%** | 95.9% | **97.0%** | 35.9% | **58.5%** | 343.73 |
| | | 80% | **29.19** | **89.32** | **9.19** | **66.2%** | 95.8% | **97.0%** | 35.7% | 58.2% | 295.31 |
| | | 90% | **24.22** | **79.90** | **18.12** | 65.9% | 95.7% | 96.9% | 36.6% | 57.8% | 283.20 |



Figure 4: Latency results tested in 768p with H100 GPU for different sparsity ratios in attention.

## 4.2 Main Results

**Higher Generation Quality.** We provide the main results compare with the SVG method in Table 1. To perform a comprehensive study, different sparsity ratios for the attention mechanism are evaluated under various resolutions with multiple video generation model architectures. With the Wan2.1 model, we observe that our method achieves less image quality degradation compared with SVG. The similarity results measured by PSNR, SSIM and LPIPS demonstrate that our method generates videos more similar to the dense model compared with SVG under the same sparsity. Specifically, for Wan2.1 (768p), our method achieves non-marginal improvements over SVG on PSNR, SSIM and LPIPS (such as our 8.71 LPIPS *v.s.* 17.57 LPIPS from SVG under 75% sparsity). For the Hunyuan model, our method achieves better performance across almost all reported metrics, under a fair comparison with SVG following the same sparsity and computational cost in PFLOPs. Although SVG includes additional overhead for spatial or temporal head selection, we exclude this computation cost from the reported PFLOPSs of SVG in Table 1. Note that the additional overhead of our `DraftAttention` is minor, leading to almost the same computations as SVG in the table.

**Superior Inference Acceleration.** Furthermore, we provide our latency results in Figure 4. The latency results are tested on H100 for both Huyuan and Wan2.1 models in 768p resolution. Our method achieves over $1.75\times$ acceleration on an H100 GPU with 90% sparsity in the attention mechanism—demonstrating our outstanding practical efficiency.

**Better Visualization.** We provide the visualization for the comparison between `DraftAttention` and SVG in Figure 5. All videos are generated with 90% sparsity in sparse attention. As highlighted in the red box, SVG exhibits a noticeable degradation in generation quality, with apparent blurry pixels. In contrast, our method better maintains the generation quality with videos more similar to the dense baseline. We provide generated videos for further visualization comparison in supplementary.
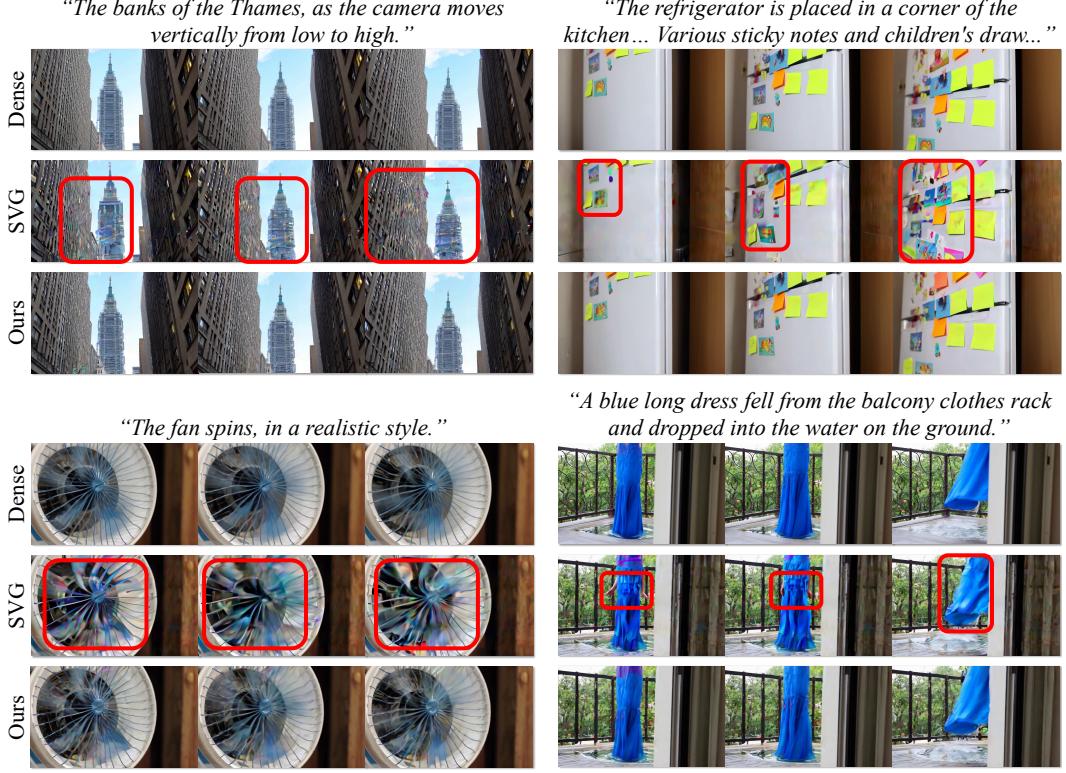
Figure 5: Visualization for our method and SVG [16] with 90% sparsity ratio in attention.
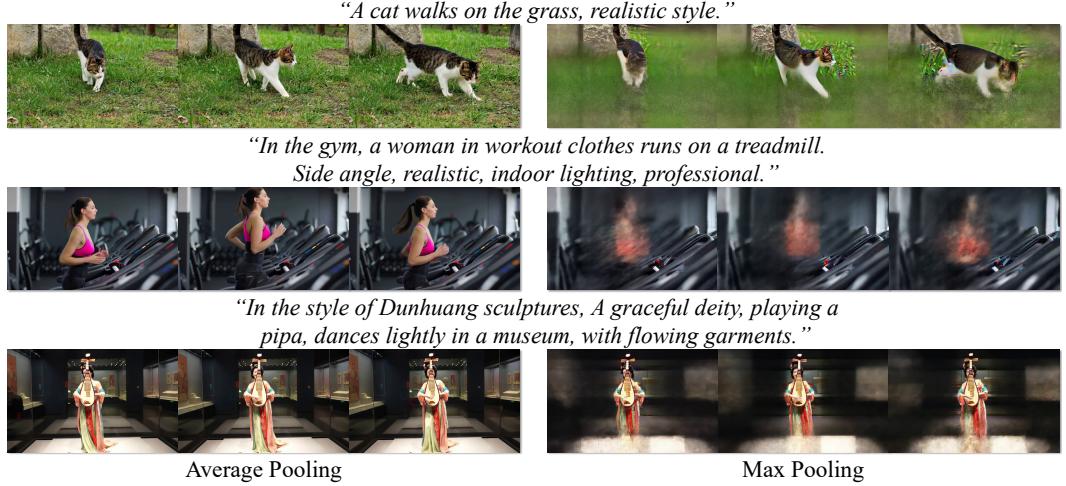


Figure 6: Visualization from the ablation study comparing average pooling and max pooling kernels for downsampling in the draft attention module with 90% sparsity.

## 4.3 Ablation Study

As shown in Figure 6, we provide the visualization of ablation study for the different downsampling kernels with average pooling and max pooling. The visualization is generated using 90% sparsity in the sparse attention, with only the average pooling replaced by max pooling in our framework. We observe that the average pooling achieves much better generation quality than the max pooling, especially for the background part.

## 5 Conclusion

In this paper, we propose the `DraftAttention` for the fast video diffusion. We adopt pooling to compute a low-resolution draft attention map to guide the sparse attention over full-resolution query, key, and value representations. Combined with effective reordering, this approach achieves fast, hardware-friendly execution on GPUs. Theoretical analysis is further provided for the justification of our design. Experiments show that our method outperforms other methods and achieves up to $1.75\times$ end-to-end acceleration on GPUs. In the future work, we plan to introduce the quantization for the further acceleration of high-resolution and long-duration video generation on GPUs.

## References

[1] William Peebles and Saining Xie. Scalable diffusion models with transformers. *arXiv preprint arXiv:2212.09748*, 2022.

[2] 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.

[3] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6836–6846, 2021.

[4] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, , et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024.

[5] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, et al. Hunyuanvideo: A systematic framework for large video generative models, 2024.

[6] Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, Jinkai Wang, Jixuan Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Mengyang Feng, Ningyi Zhang, Pandeng Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei Zhang, Siyang Sun, Tao Fang, Tianxing Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang, Wenmeng Zhou, Wente Wang, Wenting Shen, Wenyuan Yu, Xianzhong Shi, Xiaoming Huang, Xin Xu, Yan Kou, Yangyu Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yitong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi Jiang, Zhen Han, Zhi-Fan Wu, and Ziyu Liu. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025.

[7] Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.

[8] Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuandong Tian, Christopher Re, Clark Barrett, Zhangyang Wang, and Beidi Chen. H2o: Heavy-hitter oracle for efficient generative inference of large language models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.

[9] Jiaming Tang, Yilong Zhao, Kan Zhu, Guangxuan Xiao, Baris Kasikci, and Song Han. Quest: Query-aware sparsity for efficient long-context llm inference, 2024.

[10] Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models with attention sinks. *arXiv*, 2023.

[11] Huiqiang Jiang, Yucheng Li, Chengruidong Zhang, Qianhui Wu, Xufang Luo, Surin Ahn, Zhenhua Han, Amir H. Abdi, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. MInference 1.0: Accelerating pre-filling for long-context LLMs via dynamic sparse attention. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.

[12] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.

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

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

[15] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

[16] Haocheng Xi, Shuo Yang, Yilong Zhao, Chenfeng Xu, Muyang Li, Xiuyu Li, Yujun Lin, Han Cai, Jintao Zhang, Dacheng Li, et al. Sparse videogen: Accelerating video diffusion transformers with spatial-temporal sparsity. *arXiv preprint arXiv:2502.01776*, 2025.

[17] Yifei Xia, Suhan Ling, Fangcheng Fu, Yujie Wang, Huixia Li, Xuefeng Xiao, and Bin Cui. Training-free and adaptive sparse attention for efficient long video generation. *arXiv preprint arXiv:2502.21079*, 2025.

[18] Junxian Guo, Haotian Tang, Shang Yang, Zhekai Zhang, Zhijian Liu, and Song Han. Block Sparse Attention. `https://github.com/mit-han-lab/Block-Sparse-Attention`, 2024.

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

[20] Jintao Zhang, Jia wei, Pengle Zhang, Jun Zhu, and Jianfei Chen. Sageattention: Accurate 8-bit attention for plug-and-play inference acceleration. In *The Thirteenth International Conference on Learning Representations*, 2025.

[21] Jintao Zhang, Haofeng Huang, Pengle Zhang, Jia Wei, Jun Zhu, and Jianfei Chen. Sageattention2: Efficient attention with thorough outlier smoothing and per-thread int4 quantization. In *International Conference on Machine Learning (ICML)*, 2025.

[22] Muyang Li*, Yujun Lin*, Zhekai Zhang*, Tianle Cai, Xiuyu Li, Junxian Guo, Enze Xie, Chenlin Meng, Jun-Yan Zhu, and Song Han. Svdquant: Absorbing outliers by low-rank components for 4-bit diffusion models. In *The Thirteenth International Conference on Learning Representations*, 2025.

[23] Enze Xie, Junsong Chen, Junyu Chen, Han Cai, Haotian Tang, Yujun Lin, Zhekai Zhang, Muyang Li, Ligeng Zhu, Yao Lu, and Song Han. SANA: Efficient high-resolution text-to-image synthesis with linear diffusion transformers. In *The Thirteenth International Conference on Learning Representations*, 2025.

[24] Junyu Chen, Han Cai, Junsong Chen, Enze Xie, Shang Yang, Haotian Tang, Muyang Li, and Song Han. Deep compression autoencoder for efficient high-resolution diffusion models. In *The Thirteenth International Conference on Learning Representations*, 2025.

[25] 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 *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.

[26] Tianwei Yin, Michaël Gharbi, Taesung Park, Richard Zhang, Eli Shechtman, Fredo Durand, and William T. Freeman. Improved distribution matching distillation for fast image synthesis. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.

[27] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.

[28] Guangxuan Xiao, Jiaming Tang, Jingwei Zuo, junxian guo, Shang Yang, Haotian Tang, Yao Fu, and Song Han. Duoattention: Efficient long-context LLM inference with retrieval and streaming heads. In *The Thirteenth International Conference on Learning Representations*, 2025.

[29] Ruyi Xu, Guangxuan Xiao, Haofeng Huang, Junxian Guo, and Song Han. Xattention: Block sparse attention with antidiagonal scoring. *arXiv preprint arXiv:2503.16428*, 2025.

[30] Muyang Li, Tianle Cai, Jiaxin Cao, Qinsheng Zhang, Han Cai, Junjie Bai, Yangqing Jia, Ming-Yu Liu, Kai Li, and Song Han. Distrifusion: Distributed parallel inference for high-resolution diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.

[31] Muyang Li, Yujun Lin, Zhekai Zhang, Tianle Cai, Xiuyu Li, Junxian Guo, Enze Xie, Chenlin Meng, Jun-Yan Zhu, and Song Han. Svdqunat: Absorbing outliers by low-rank components for 4-bit diffusion models. *arXiv preprint arXiv:2411.05007*, 2024.

[32] 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.

[33] 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.

[34] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.

# Appendix

## A Detailed Proof

### A.1 Proof of Theorem 3.3

*Proof.* First, observe that for any $u \in R_i$ and $v \in R_j$, the draft attention assigns

$$(S_{\mathrm{draft}})_{uv} = \widetilde{S}_{ij}, \quad \text{while} \quad S_{uv} = \langle Q_u, K_v \rangle. \tag{12}$$

By the definition of $\delta$, we have

$$|S_{uv} - (S_{\mathrm{draft}})_{uv}| = |S_{uv} - \widetilde{S}_{ij}| \leq \delta. \tag{13}$$

Then, summing over all $n^2$ token pairs gives

$$\|S - S_{\mathrm{draft}}\|_F^2 = \sum_{u,v} |S_{uv} - (S_{\mathrm{draft}})_{uv}|^2 \leq n^2 \delta^2. \tag{14}$$

Taking square roots on both sides yields the desired result:

$$\|S - S_{\mathrm{draft}}\|_F \leq \delta n. \tag{15}$$

This completes the proof. $\square$

### A.2 Proof of Theorem 3.5

*Proof.* The mask $\widehat{M}$ zeros out exactly $(1 - r)n^2$ entries corresponding to dropped blocks. For each $(u, v)$ in a dropped block, we have:

$$|S_{uv}| \leq |S_{uv} - \widetilde{S}_{ij}| + |\widetilde{S}_{ij}| \leq \delta + t. \tag{16}$$

Summing squared errors over $(1 - r)n^2$ entries yields the bound:

$$\|S - S \odot \widehat{M}\|_F^2 \leq (1 - r)n^2(\delta + t)^2 \quad \Rightarrow \quad \|S - S \odot \widehat{M}\|_F \leq n(\delta + t)\sqrt{1 - r}. \tag{17}$$

We then finish the proof. $\square$