

**KAIST**

 DATA MINING LAB

**CVPR**  
JUNE 3-7, 2026



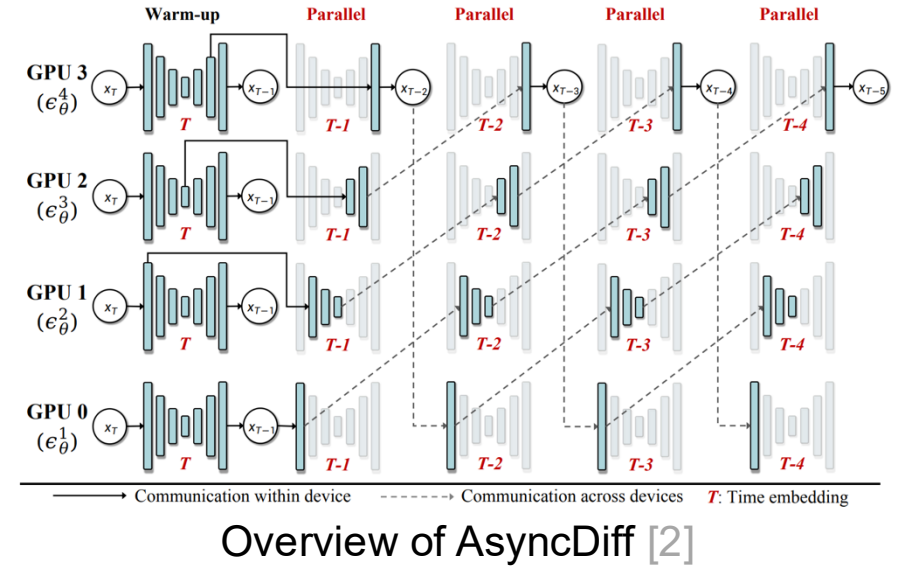
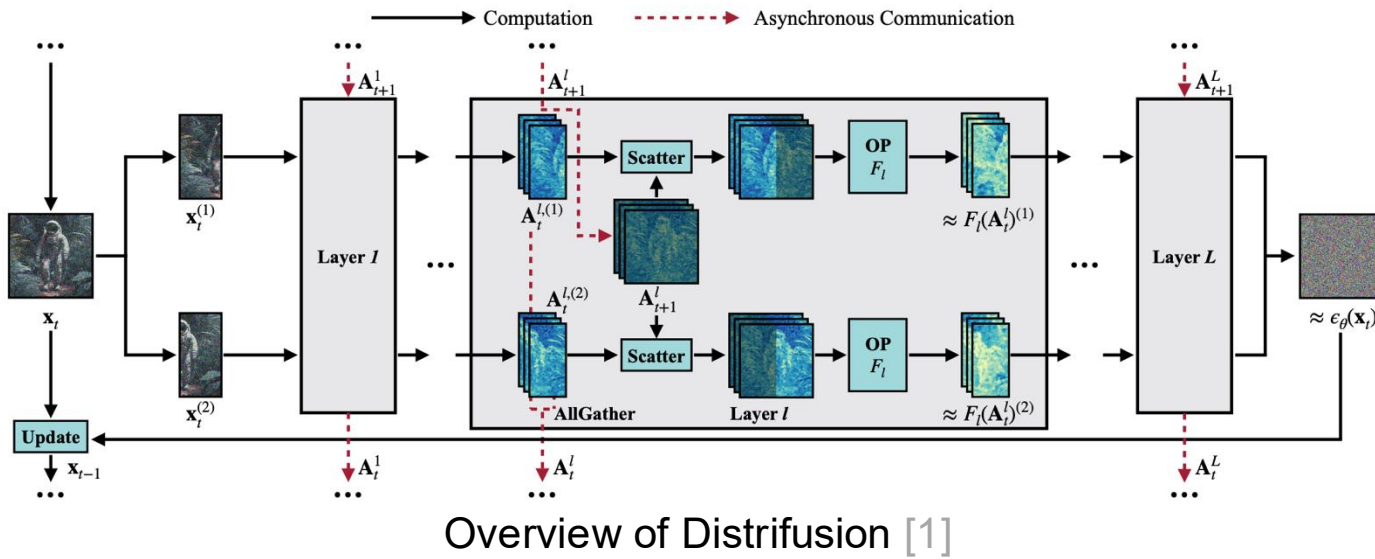
**DENVER**  
COLORADO

# **Accelerating Diffusion via Hybrid Data-Pipeline Parallelism Based on Conditional Guidance Scheduling**

CVPR 2026 poster

Euisoo Jung, Byunghyun Kim, Hyunjin Kim, Seonghye Cho, Jae-Gil Lee

# Accelerating Diffusion Inference via Multi-GPU

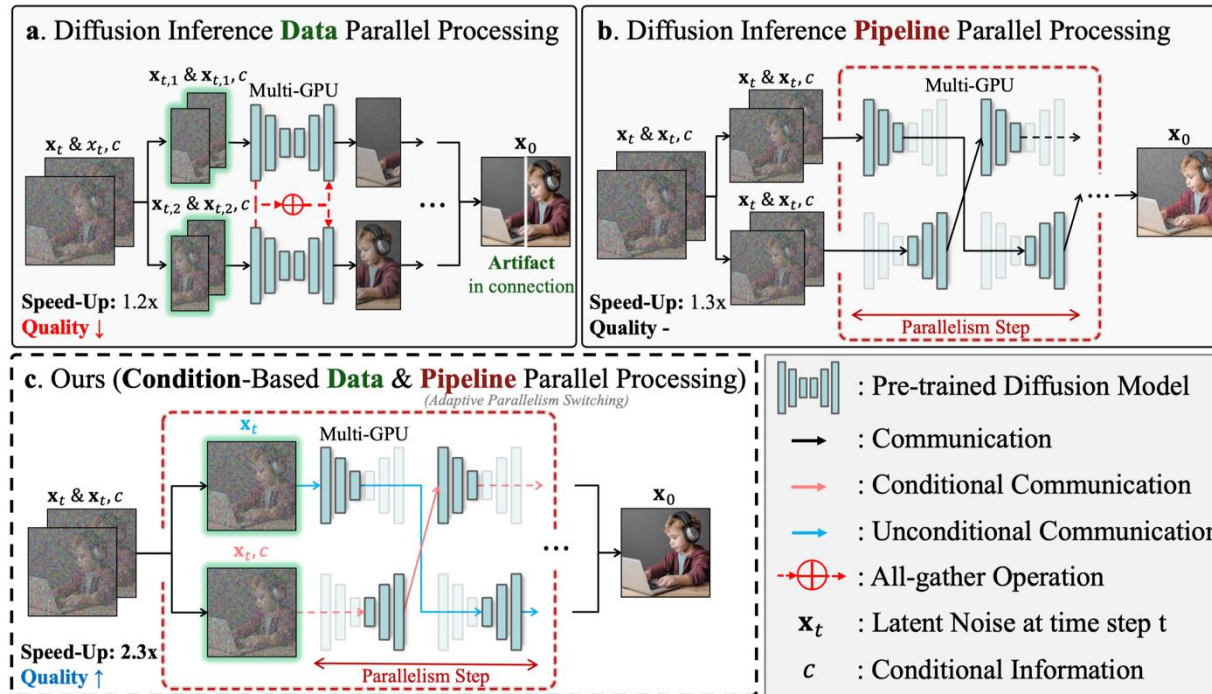


- **Diffusion models** have become the backbone of image and video generation, but their **long inference latency** remains a key bottleneck, motivating extensive acceleration research
- Among these, **multi-GPU distributed inference** has emerged as a promising direction for accelerating generation
- Representative approaches include **Distrifusion** (patch-based data parallelism combined with tensor parallelism) and **AsyncDiff** (pipeline parallelism via asynchronous denoising)

[1] Li, Muyang, et al. "Distrifusion: Distributed parallel inference for high-resolution diffusion models." *CVPR*. 2024.

[2] Chen, Z., Ma, X., Fang, G., Tan, Z., & Wang, X. Asyncdiff: Parallelizing diffusion models by asynchronous denoising. *NeurIPS*. 2024.

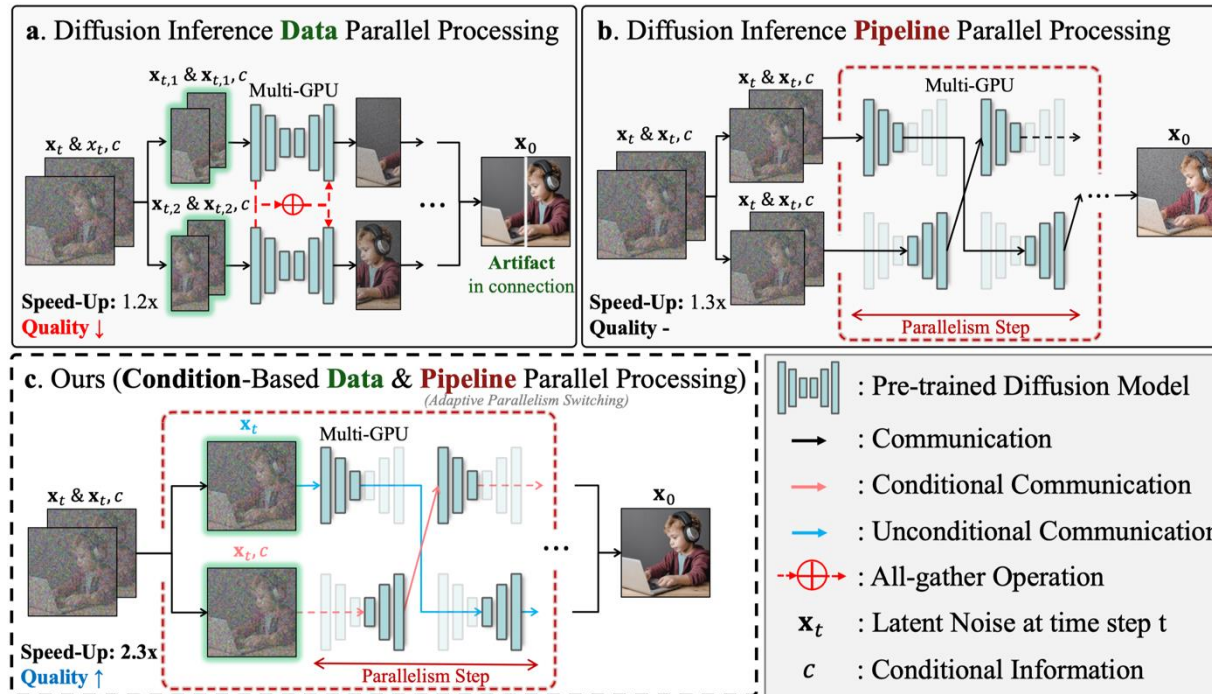
# Prior Multi-GPU Inference Works Problems



Comparison of parallel strategies for diffusion inference

- Multi-GPU distributed inference has established itself as a **training-free yet powerful acceleration** approach, but existing methods still suffer from several limitations
  - **Patch-based data parallelism** can degrade output quality due to **artifacts at patch connection boundaries**
  - **Pipeline parallelism** suffers from error accumulation caused by **heavy asynchronous communication**

# Prior Multi-GPU Inference Works Problems

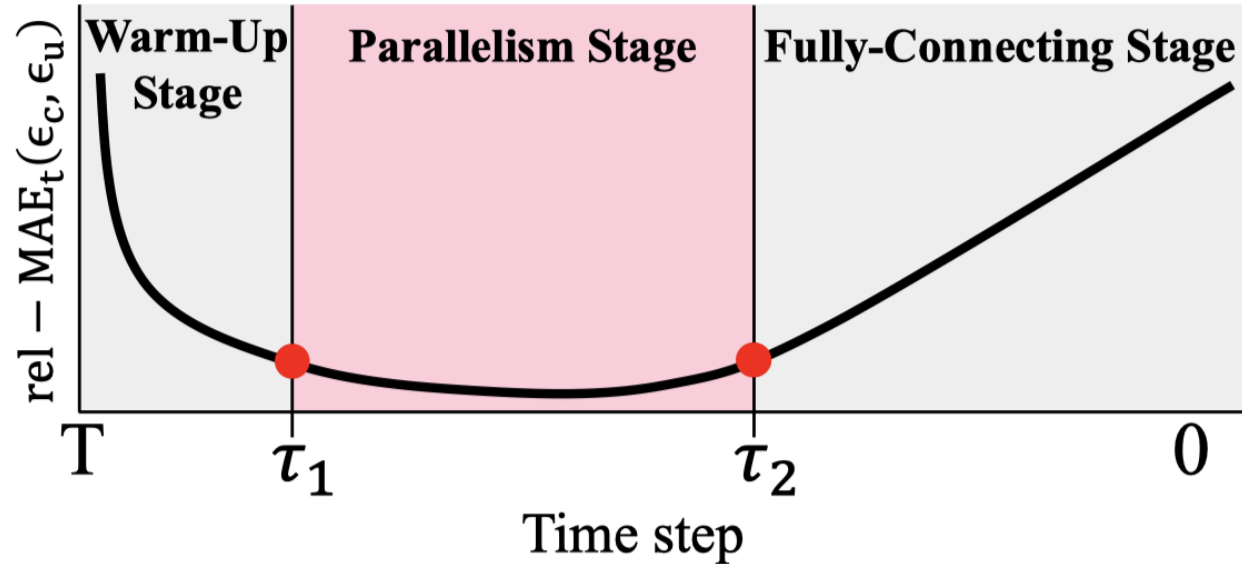


Comparison of parallel strategies for diffusion inference

✅ **Key Insight (TL;DR)** : The conditional & unconditional CFG paths form a natural, artifact-free data split for data parallelism;

combined with pipeline parallelism in the interval where their discrepancy is minimal, they constitute a hybrid parallelism that achieves beyond-linear 2.3× speed-up on 2 GPUs

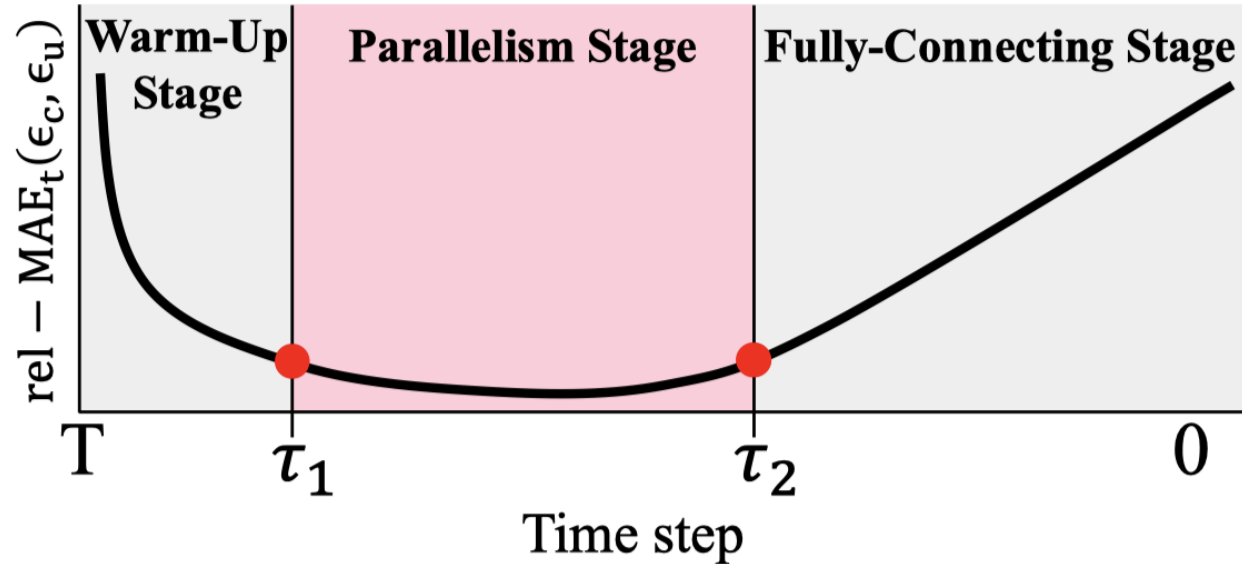
# Classifier-Free Guidance (CFG) Analysis



relative -  $\text{MAE}_t(\epsilon_c, \epsilon_u)$  curve according to denoising step

- **CFG predicts two noises** at every step: **conditional** ( $\epsilon_c$ ) and **unconditional** ( $\epsilon_u$ )  
We use their gap to decide when parallelism is safe
- **Denoising discrepancy**:  $\text{relative} - \text{MAE}_t(\epsilon_c, \epsilon_u) = \mathbb{E}[\|\epsilon_c - \epsilon_u\|_1] / \mathbb{E}[\|\epsilon_u\|_1]$   
→ How far the two paths diverge at step  $t$
- Over 5,000 MS-COCO prompts, this curve forms a U-shape, splitting denoising into three stages:
  - **Warm-Up** [ $T, \tau_1$ ]: **large gap** → branches run independently
  - **Parallelism** ( $\tau_1, \tau_2$ ): **small / stable gap** → **parallelism activated** for max speed-up (branch merged)
  - **Fully-Connecting** [ $\tau_2, 0$ ]: **gap rises again** → branches split for refinement

# Classifier-Free Guidance (CFG) Analysis

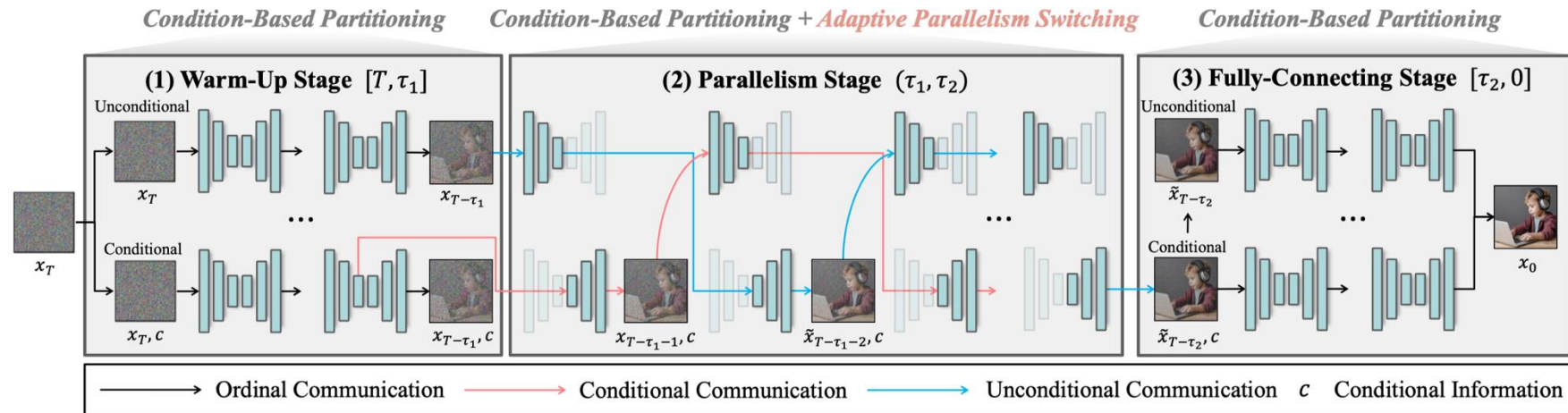


relative -  $\text{MAE}_t(\epsilon_c, \epsilon_u)$  curve according to denoising step

- Over 5,000 MS-COCO prompts, this curve forms a U-shape, splitting denoising into three stages:
  - **Warm-Up** [ $T, \tau_1$ ]: **large gap** → branches run independently
  - **Parallelism** ( $\tau_1, \tau_2$ ): **small / stable gap** → **parallelism activated** for max speed-up (branch merged)
  - **Fully-Connecting** [ $\tau_2, 0$ ]: **gap rises again** → branches split for refinement

➔ Need to find  $\tau_1, \tau_2$  dynamic & automatically

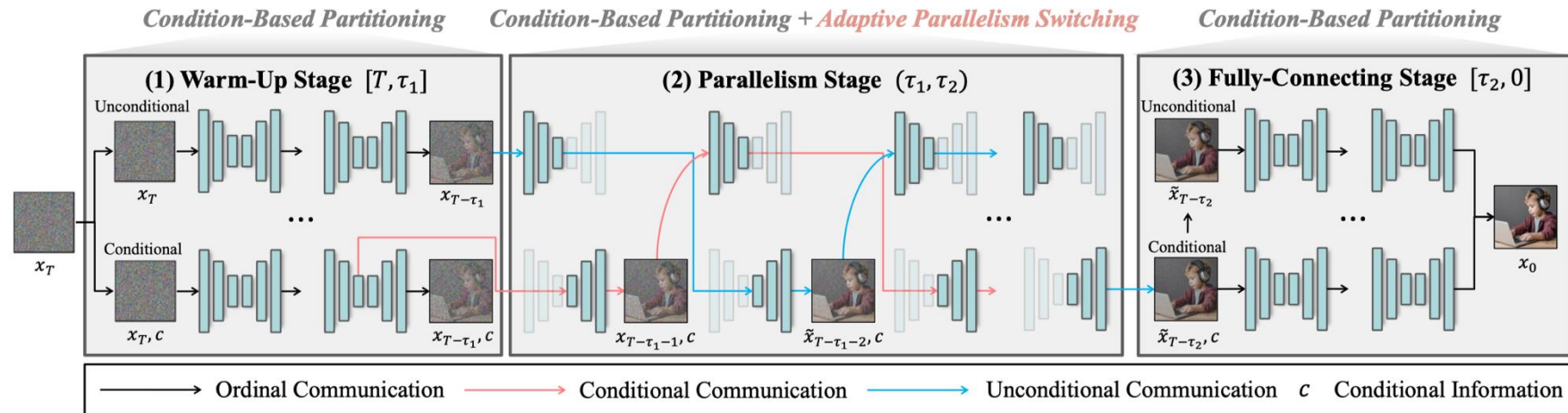
# Method: Hybrid Data-Pipeline Parallelism



Overview of proposed hybrid data-pipeline parallelism

- *Condition-Based Partitioning*
  - Split by CFG branch (conditional & unconditional)  
→ Enable data parallelism **w/o artifacts**
- *Adaptive Parallelism Switching*
  - Engage pipeline parallelism only in  $[\tau_1, \tau_2]$  where two branches similar  
→ Enable **hybrid data-pipeline parallelism** w/o minimize the quality degradation

# Method: Hybrid Data-Pipeline Parallelism



Overview of proposed hybrid data-pipeline parallelism

- Determining  $\tau_1$ 
  1.  $M_t = \text{relative} - \text{MAE}_t(\epsilon_c, \epsilon_u)$
  2.  $G_t = \frac{M_t - M_{t-L}}{L}$
  3.  $\tau'_1 = \min\{t \mid 0 \leq G_t < g_{\text{slope}}\}$
  4.  $\tau_1 = \min(\tau'_1, \tau_{\text{cap}})$
- Determining  $\tau_2$ 
  1.  $\tau_2 = \tau_1 + k, \quad k \in \mathbb{N}, 1 \leq k < T - \tau_1$

Detailed  
Algorithm



## Algorithm 1 Adaptive Parallelism Switching via Denoising Discrepancy

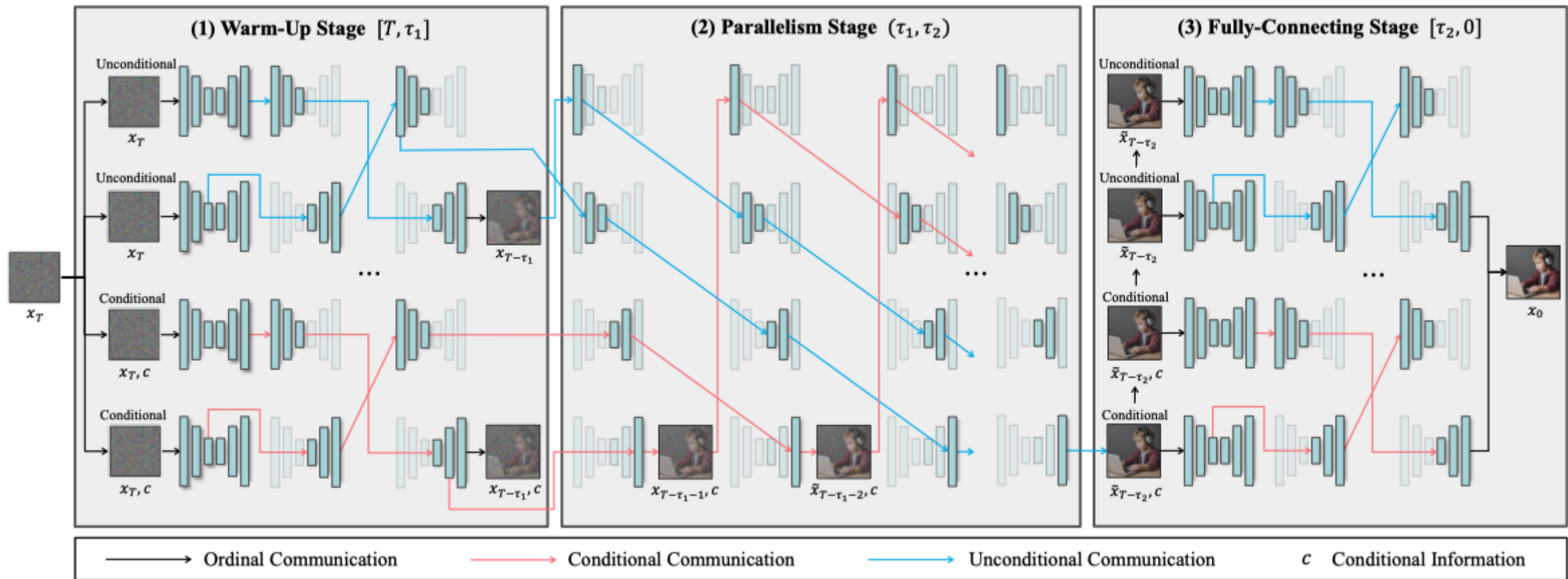
**Require:** latent noise  $\mathbf{x}_t$ , prompt  $c$ , steps  $T$ , window  $L$ , slope threshold  $g$ , safety-cap  $\tau_{\text{cap}}$ , interval  $k$

```

1:  $\tau_1, \tau_2 \leftarrow \emptyset$ 
2: for  $t = T, T-1, \dots, 1$  do
3:    $\epsilon_c, \epsilon_u \leftarrow \epsilon_\theta(\mathbf{x}_t, c, t), \epsilon_\theta(\mathbf{x}_t, t)$ 
4:    $M_t \leftarrow \frac{\mathbb{E}_{\mathbf{x}, \epsilon} \|\epsilon_c - \epsilon_u\|_1}{\mathbb{E}_{\mathbf{x}, \epsilon} \|\epsilon_u\|_1} \triangleright \text{rel-MAE}_t(\epsilon_c, \epsilon_u)$ 
5:    $G_t = \frac{M_t - M_{t-L}}{L}$ 
6:   if  $\tau_1 = \emptyset$  and  $t > L$  and  $0 \leq G_t < g$  then
7:      $\tau_1 \leftarrow \min(t, \tau_{\text{cap}}); \tau_2 \leftarrow \tau_1 + k$ 
8:   Denoise:
9:   if  $t \geq \tau_1$  then
10:     WARM-UP
11:   else if  $t > \tau_2$  then
12:     PARALLELISM
13:   else
14:     FULLY-CONNECTING
15:   end if
16:    $x_{t-1} \leftarrow \text{STEP DENOISE}(\mathbf{x}_t, \epsilon_c, \epsilon_u, t)$ 
17: end for
18: return  $x_0, (\tau_1, \tau_2)$ 

```

# Method: Hybrid Data-Pipeline Parallelism



Overview of proposed hybrid data-pipeline parallelism in over 2 GPUs

# Main Results: Quantitative Results

Base Model	Devices	Methods	Latency (s) ↓	Speed-Up ↑	Comm. (GB) ↓	FID ↓		LPIPS ↓		PSNR ↑		
						w/ G. T.	w/ Orig.	w/ G. T.	w/ Orig.	w/ G. T.	w/ Orig.	
Stable Diffusion XL	1	Original Model	16.49	-	-	23.977	-	0.797	-	9.618	-	
	2	DistriFusion [12]	13.53	1.22×	0.525	24.164	4.864	0.7978	0.146	9.597	24.634	
		AsyncDiff [2] (stride=1)	12.54	1.31×	9.830	23.941	4.103	0.797	0.108	9.586	26.387	
		<b>Ours (<math>k=5</math>)</b>	<b>7.12</b>	<b>2.31×</b>	<b>0.516</b>	<b>23.831</b>	<b>4.100</b>	<b>0.796</b>	<b>0.107</b>	<b>9.665</b>	<b>26.640</b>	
	4	DistriFusion	7.13	2.31×	<b>0.486</b>	24.197	5.772	0.798	0.183	9.578	23.046	
		AsyncDiff (stride=1)	9.45	1.74×	10.531	24.225	5.829	0.795	0.147	9.626	24.776	
		AsyncDiff (stride=2)	6.57	2.51×	4.795	24.140	6.572	0.803	0.192	9.557	23.070	
		<b>Ours (<math>k=5</math>)</b>	<b>4.83</b>	<b>3.41×</b>	<b>0.751</b>	<b>24.087</b>	<b>5.544</b>	<b>0.788</b>	<b>0.113</b>	<b>9.888</b>	<b>25.233</b>	
	Stable Diffusion 3	1	Original Model	19.36	-	-	33.433	-	0.810	-	8.086	-
		2	AsyncDiff [2] (stride=1)	9.82	1.97×	1.290	33.379	2.032	0.813	0.052	8.155	27.812
xDiT-Ring [4]			14.31	1.35×	121.646	33.356	1.909	0.809	0.047	8.085	27.857	
Parastep [33]			9.98	1.94×	<b>0.032</b>	33.340	3.350	0.810	0.112	8.091	22.917	
Compact-2bit [21]			13.96	1.39×	9.049	33.560	4.712	0.810	0.196	7.999	19.242	
<b>Ours (<math>k=5</math>)</b>			<b>9.33</b>	<b>2.07×</b>	<b>0.189</b>	<b>33.322</b>	<b>1.878</b>	<b>0.780</b>	<b>0.046</b>	<b>8.229</b>	<b>27.875</b>	
4		AsyncDiff (stride=1)	6.51	2.97×	3.774	33.909	3.458	0.887	0.177	8.033	26.540	
		AsyncDiff (stride=2)	<b>4.86</b>	<b>3.98×</b>	1.241	34.434	8.939	0.895	0.272	7.956	21.087	
		xDiT-Ring	17.96	1.08×	130.792	33.374	2.232	0.861	0.127	8.001	27.101	
		Parastep	6.29	3.08×	<b>0.038</b>	33.211	2.581	0.910	0.134	7.990	27.089	
	Compact-2bit	9.25	2.09×	13.593	33.420	5.461	<b>0.811</b>	0.248	7.926	17.490		
<b>Ours (<math>k=5</math>)</b>	<b>5.53</b>	<b>3.50×</b>	<b>0.572</b>	<b>33.113</b>	<b>2.109</b>	0.857	<b>0.122</b>	<b>8.046</b>	<b>27.110</b>			

Main quantitative results

- On **2 GPUs**, Ours achieves **2.31×** (SDXL) and **2.07×** (SD3) speed-up  
→ outperforming DistriFusion and AsyncDiff while preserving FID/LPIPS/PSNR
- **Communication cost** is up to **19.6×** lower than AsyncDiff, and the gain scales to **3.41×** (SDXL) and **3.50×** (SD3) on **4 GPUs** → confirming generality across **U-Net and DiT** architectures.

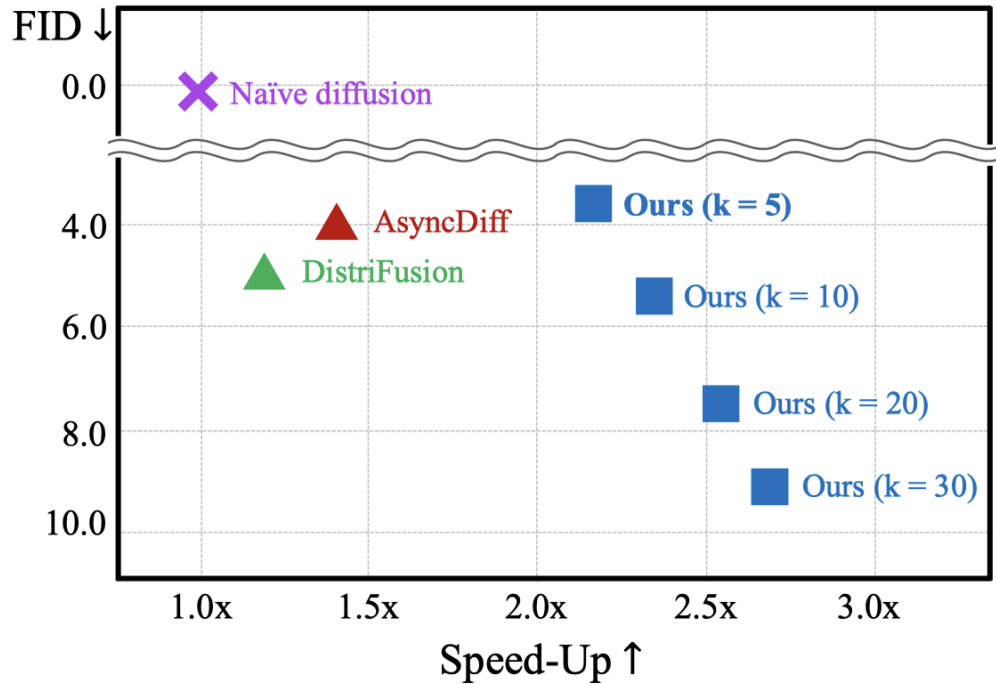
# Main Results: Qualitative Results



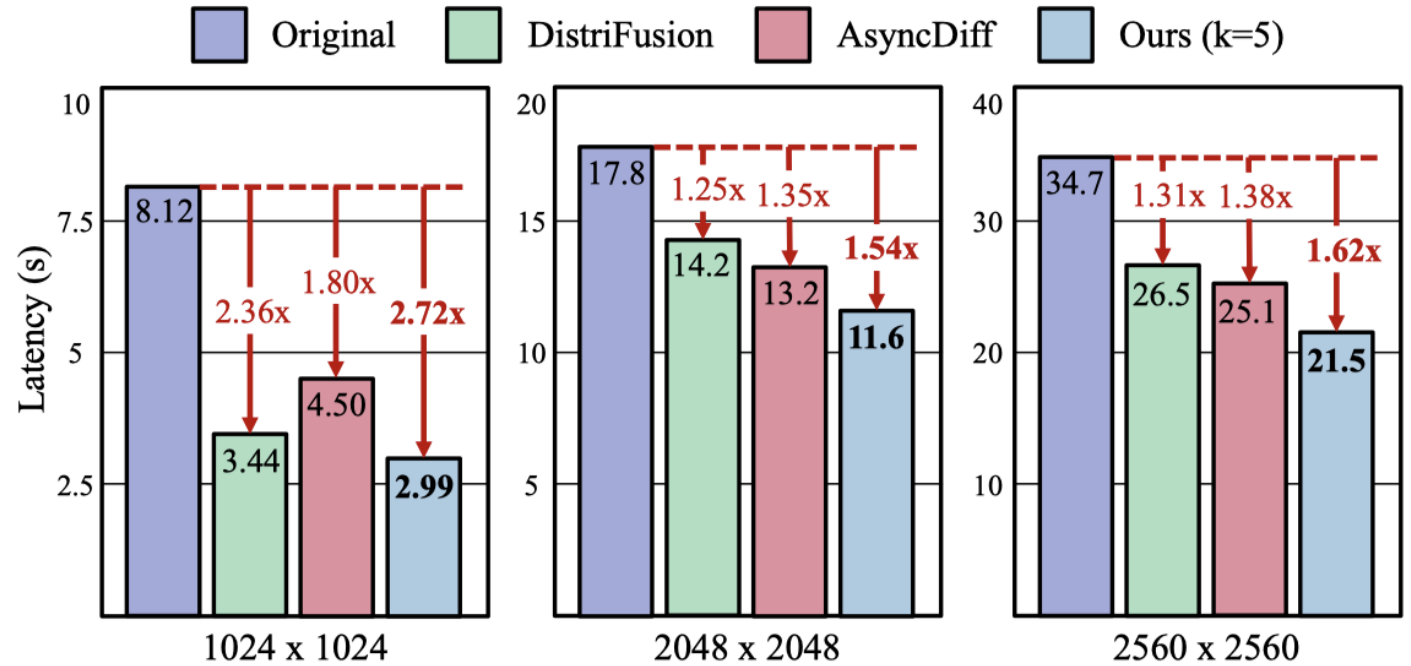
Main qualitative results

- DistriFusion shows boundary **artifacts** and AsyncDiff exhibits **spatial inconsistency**, while Ours preserves **global coherence** and **fine details** closest to the original
- Despite being the **fastest (2.32x)**, Ours also yields the **best FID (4.14)**  
→ demonstrating that condition-based partitioning avoids the quality trade-off seen in prior methods

# Sensitivity Analysis: k Values & High Resolution



Visualization of different k values



Results of high resolution tasks (H200 GPUs)

- **k controls the speed-quality trade-off:** smaller k preserves fidelity, larger k boosts speed; **k=5** sits on the Pareto frontier, dominating both prior methods
- **High-resolution scalability:** Ours maintains the best acceleration across **1024 x 1024 (2.72x)**, **2048 x 2048 (1.54x)**, and **2560 x 2560 (1.62x)** on H200 GPUs

# Summary

---

- **Method:** Hybrid parallelism combining condition-based partitioning ( $\epsilon_c, \epsilon_u$  as two data paths) with adaptive parallelism switching ( $\tau_1, \tau_2$  set by denoising discrepancy)
- **Results:**  $2.31\times$  (SDXL) and  $2.07\times$  (SD3) speed-up on 2 GPUs with preserved fidelity, scaling to  $3.41\times$  /  $3.50\times$  on 4 GPUs
- **Contribution:** A **unified, training-free** parallelism paradigm that generalizes across U-Net, DiT, and high-resolution settings with minimal communication cost.

<https://github.com/kaist-dmlab/Hybridiff>

Codes are available via online & right QR code

**Thank you!**

