

# FAST-DLLM V2: Efficient Block-Diffusion LLM

Chengyue Wu<sup>1,2</sup> Hao Zhang<sup>2</sup> Shuchen Xue<sup>2</sup> Shizhe Diao<sup>2</sup>  
Yonggan Fu<sup>2</sup> Zhijian Liu<sup>2</sup> Pavlo Molchanov<sup>2</sup> Ping Luo<sup>1</sup>  
Song Han<sup>2,3</sup> Enze Xie<sup>2</sup>

<sup>1</sup> The University of Hong Kong <sup>2</sup> NVIDIA <sup>3</sup> MIT

arXiv:2509.26328

Presented by Xiliang Xian

# Outline

---

□ **Background**

□ Design

□ Evaluations

□ Discussion

# Auto-Regressive(AR) LLM

---

## Autoregression:

Generation steps

↓  
There are three categories of the average  
There are three categories of the average rate  
↓  
There are three categories of the average rate of...

✓ High quality

✓ Arbitrary-length

✓ KV caching

✗ Not Parallelizable

□ produces one token at a time in strict left-to-right order

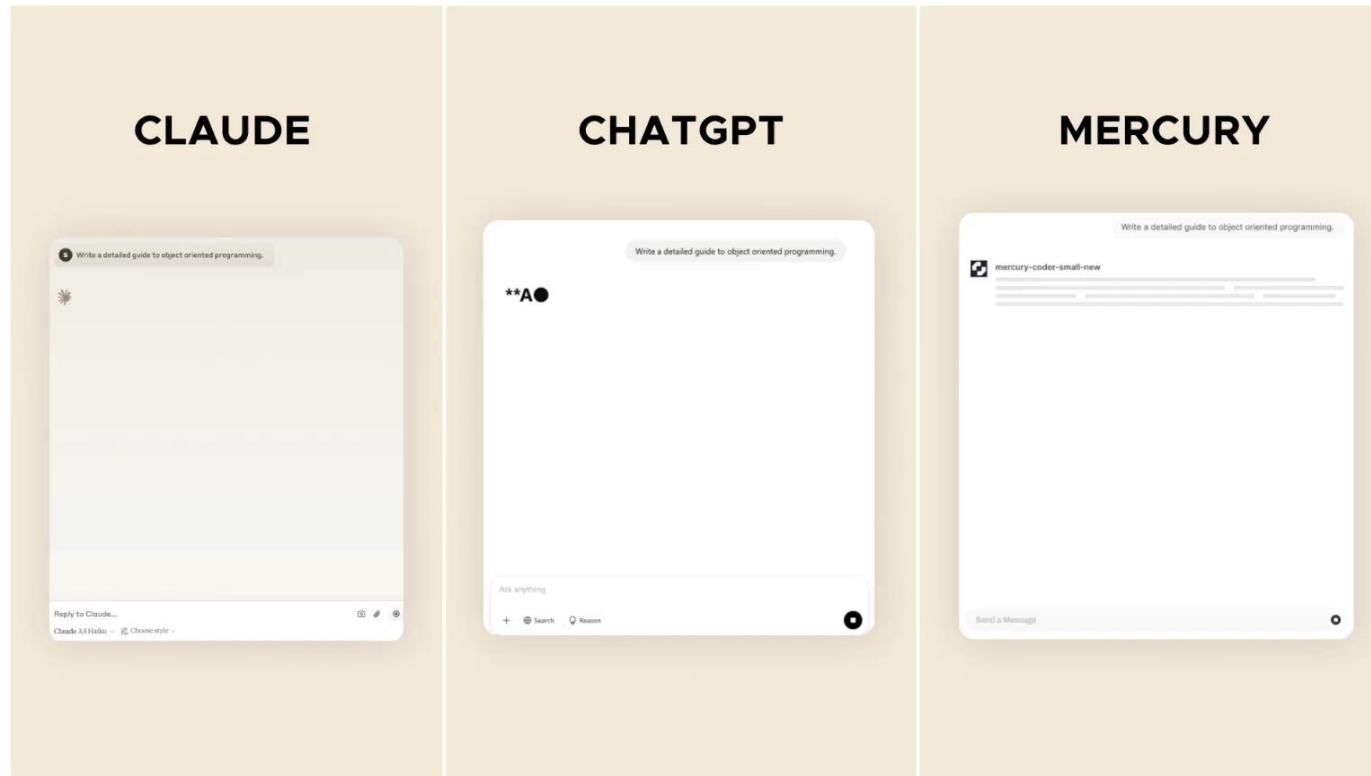
# What is DLLM?

---

- ❑ Diffusion large language model

# Why DLLM?

- ❑ DLLMs show **superior speed** vs Auto-Regressive(AR) LLM with **comparable performance**



# A Glimpse into Recent Development

---

## ❑ Closed-source model

- ❖ Gemini Diffusion
- ❖ Mercury (claimed to achieve 1109 tokens/s on H100s)
- ❖ Seed Diffusion (claimed to achieve 2146 tokens/s on H20s )

## ❑ Open-source model

- ❖ LLaDA (train from scratch)
- ❖ Dream (fine-tune from Qwen-2.5 7B)

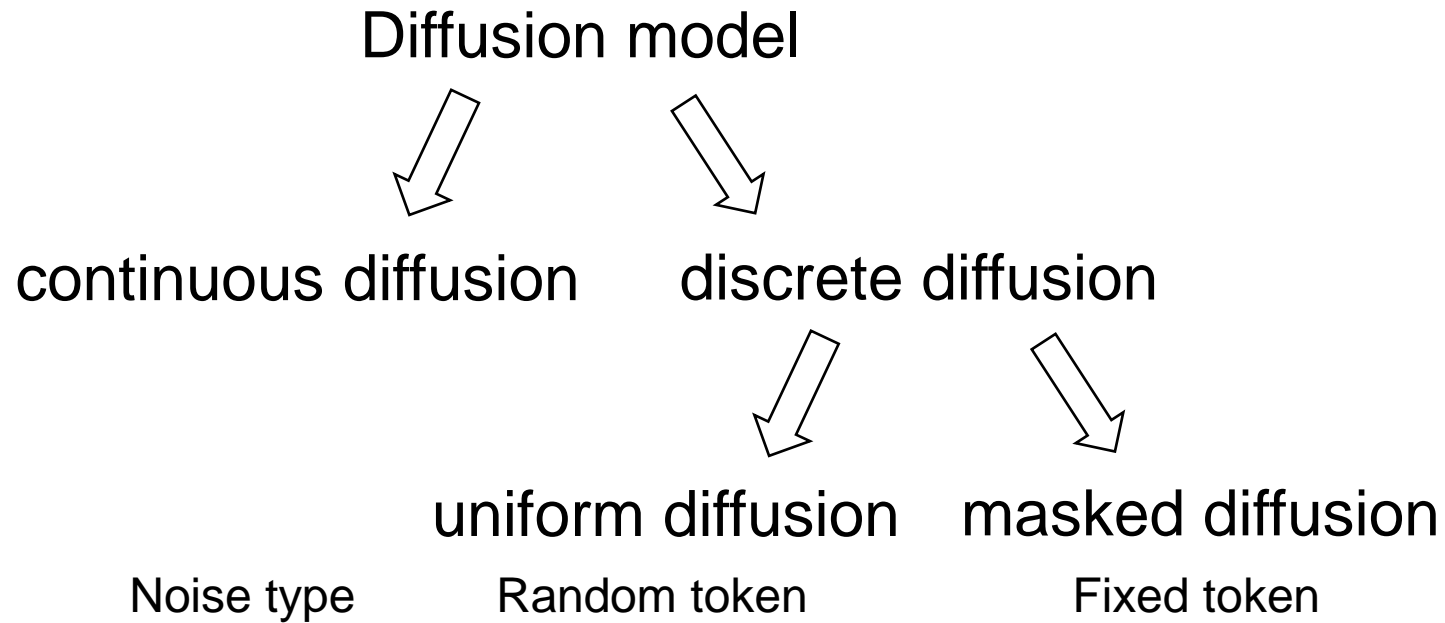
# Introduction

---

- ❑ Model Apply **Block Diffusion**
- ❑ **Data Efficient** Fine-tune Training Model
- ❑ Introduce Approximate **KV Cache** for **Full Attention**
- ❑ Introduce **Tokens** Parallel Decoding

# Classification

---





# Masked Diffusion Model (MDM)

Mainstream Formulation in DLLM

## ❑ Training details:

most Encoder-only (full attention)

❖ Forward:



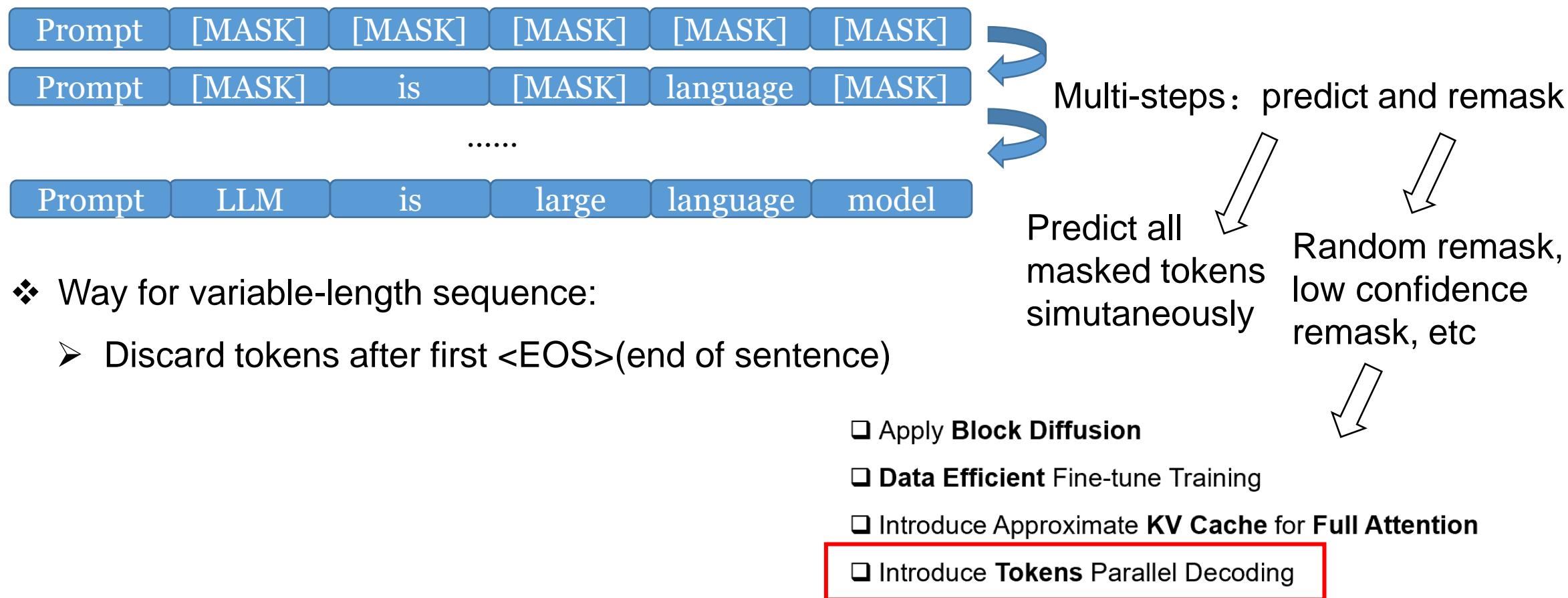
Mask Ratio  $t \sim U(0,1)$

❖ Reverse denoising:



# Masked Diffusion Model (MDM)

## ❑ Inference details:



# Block Diffusion (ICLR 25)

## ☒ Apply **Block Diffusion**

☐ **Data Efficient** Fine-tune Training

☐ Introduce Approximate **KV Cache** for **Full Attention**

☐ Introduce **Tokens** Parallel Decoding

## ☐ Semi-Autoregressive

### **Autoregression:**

✓ High quality

✓ Arbitrary-length

✓ KV caching

✗ Not Parallelizable

Generation steps

↓  
There are three categories of the average  
There are three categories of the average rate  
↓  
There are three categories of the average rate of...

### **Diffusion:**

✗ Lower quality

✗ Fixed-length

✗ No KV caching

✓ Parallelizable

↓  
the reusability will continue to the  
Repeal the reusability cuts and the law will continue to reduce the  
↓  
Repeal the reusability cuts and prove the law will continue to reduce the deficit.

### **Block Diffusion (Ours):**

✓ High quality

✓ Arbitrary-length

✓ KV caching

✓ Parallelizable

↓  
On September 17, we be  
On September 17, 2016, we will be giving the release of  
↓  
On September 17, 2016, we will be giving the beta-release of the to our server testing ...

# Limitations in MDMs

---

## ❑ Inductive bias conflict in training

- ❖ natural language is overwhelmingly processed in a sequential order

## ❑ Inference inefficiency

- ❖ take LLaDA as an example:
  - Full attention is slower than causal attention
  - **Cannot** leverage **KV cache**
  - Reduce the number of inference steps ⇒ Decode more tokens simultaneously  
⇒ **Bad accuracy**

❑ Apply **Block Diffusion**

❑ **Data Efficient** Fine-tune Training

❑ Introduce Approximate **KV Cache** for **Full Attention**

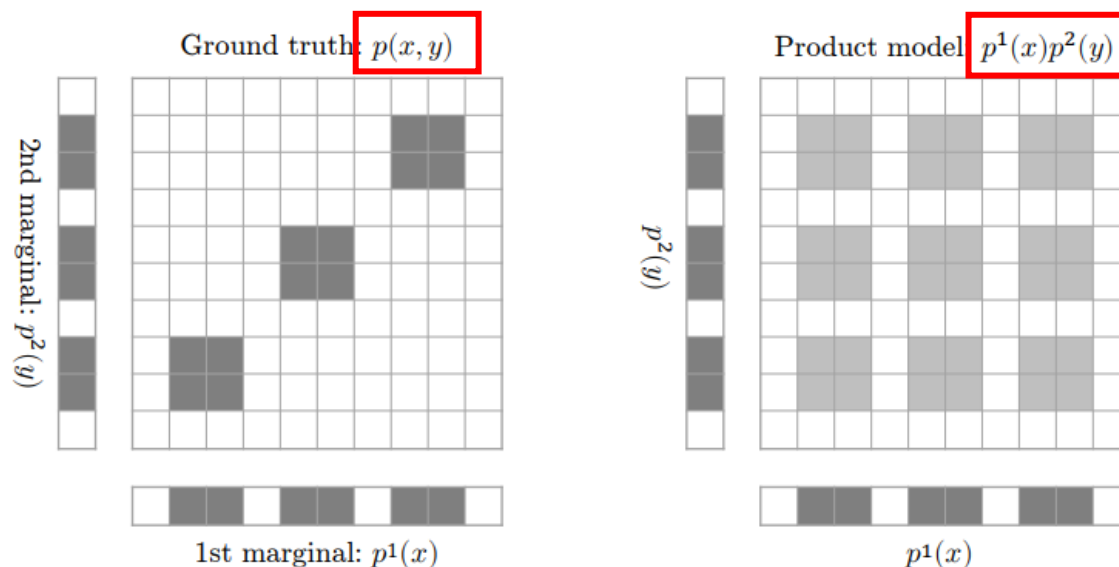
❑ Introduce **Tokens** Parallel Decoding

# Limitations in MDMs

## ❑ Inference inefficiency

.....

- Reduce the number of inference steps  $\Rightarrow$  Decode more tokens simultaneously  
 $\Rightarrow$  **Bad accuracy**



I am  
he is

**marginal** distribution  $\neq$  **joint** distribution

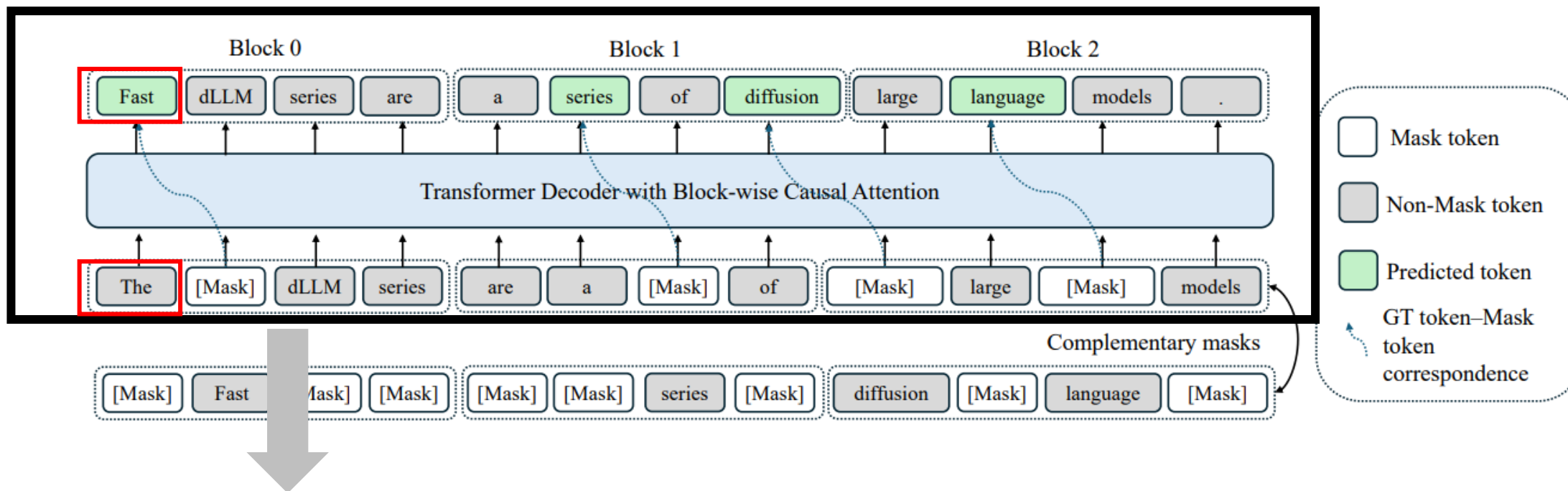
- ❑ Apply **Block Diffusion**
- ❑ **Data Efficient** Fine-tune Training
- ❑ Introduce Approximate **KV Cache** for **Full Attention**
- ❑ Introduce **Tokens Parallel Decoding**

# Outline

---

- Background
- **Design**
- Evaluations
- Discussion

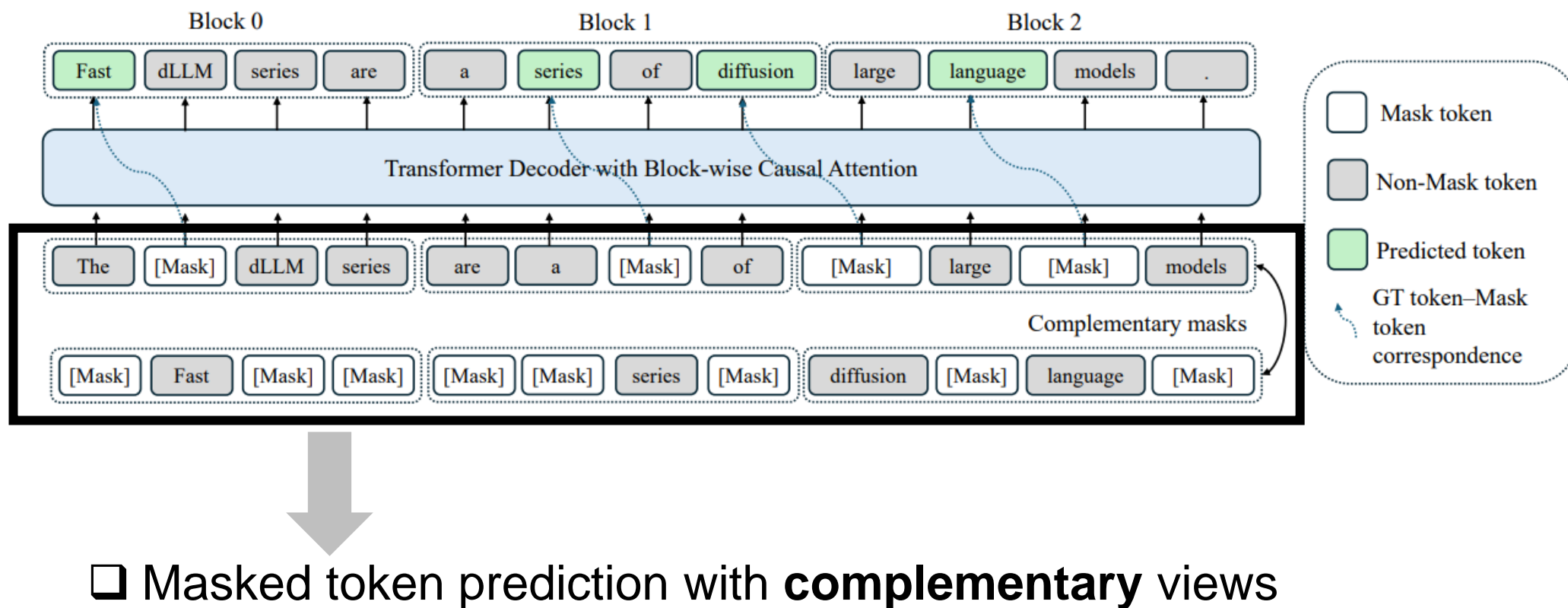
# Training Process



## Token shift for prediction

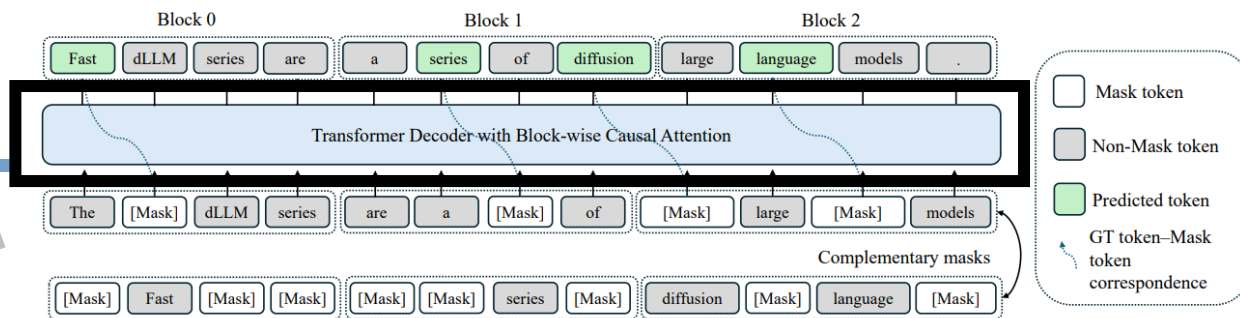
- ❖ Use hidden state at  $i-1$  to predict  $x_i$

# Training Process

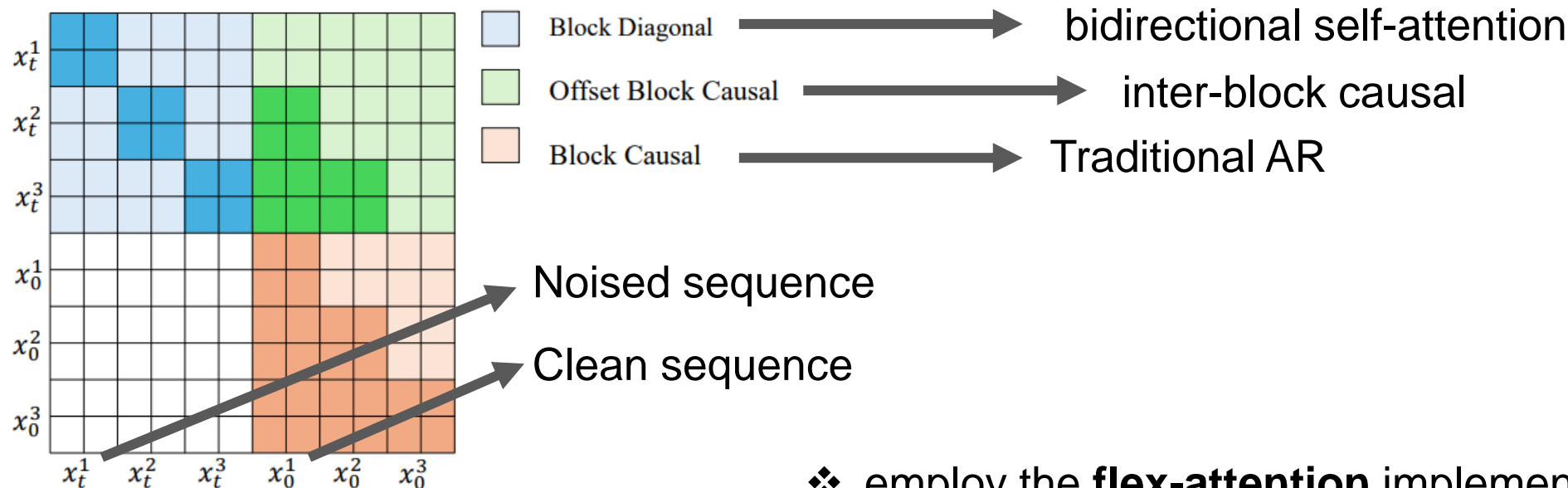




# Training Process



## Block-wise attention



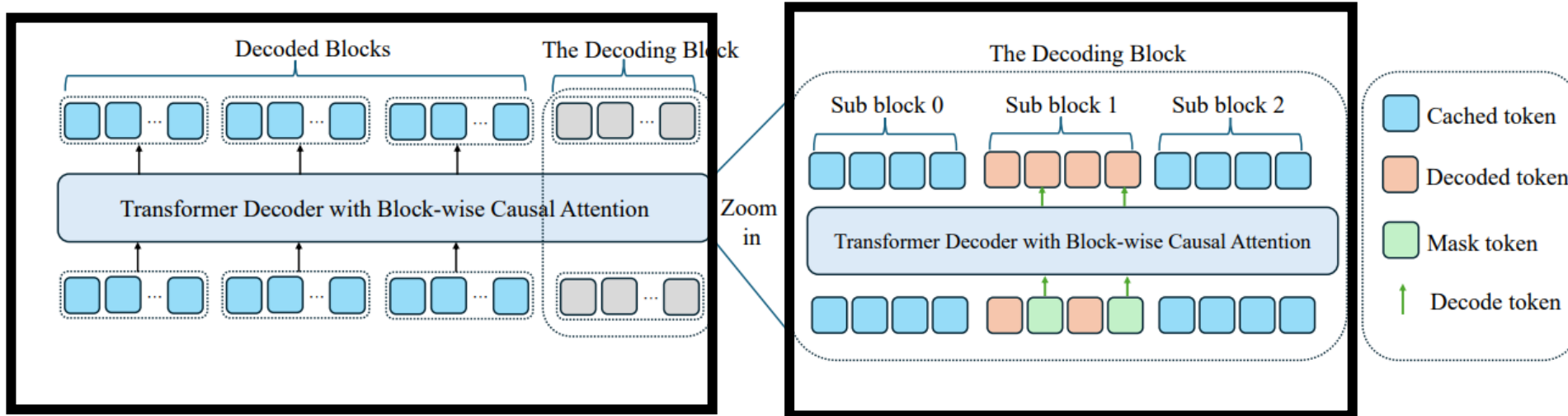
(a) Training-time attention mask.

$x^b$  denote set of tokens in the  $b$ -th **block** (rather than the  $b$ -th **token**)

❖ employ the **flex-attention** implementation

# Inference Pipeline

## ❑ Block-wise AR decoding with caching (Semi-AR)



❑ AR across blocks

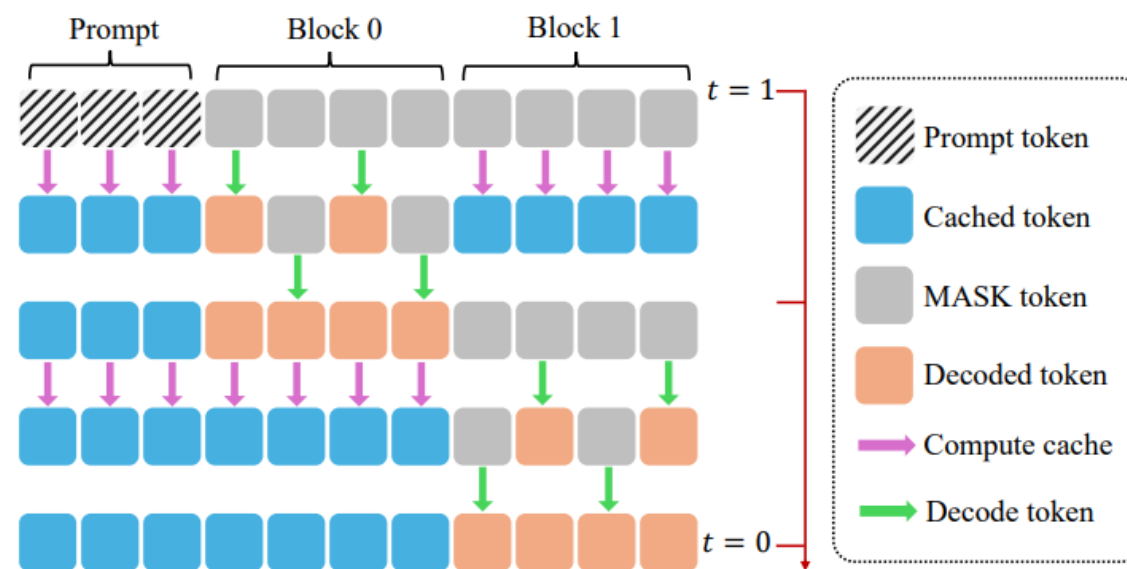
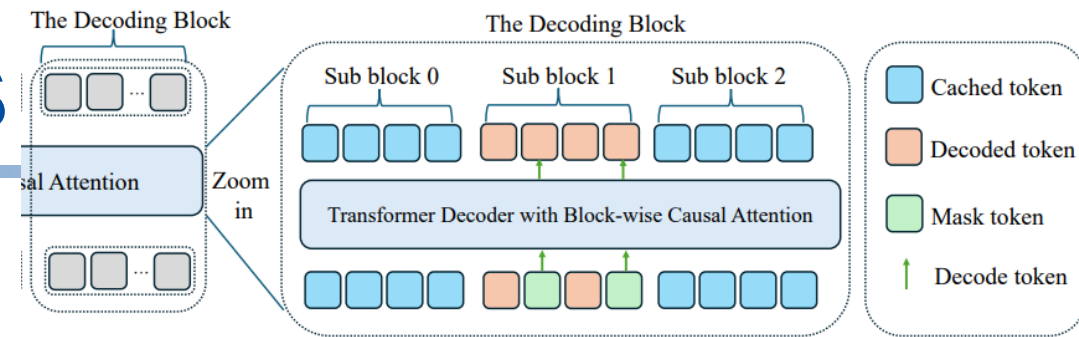
❖ Enable block-level KV cache

❑ Diffusion within blocks

# Diffusion within blocks

Cannot leverage KV cache?

❑ **DualCache**: approximate KV cache

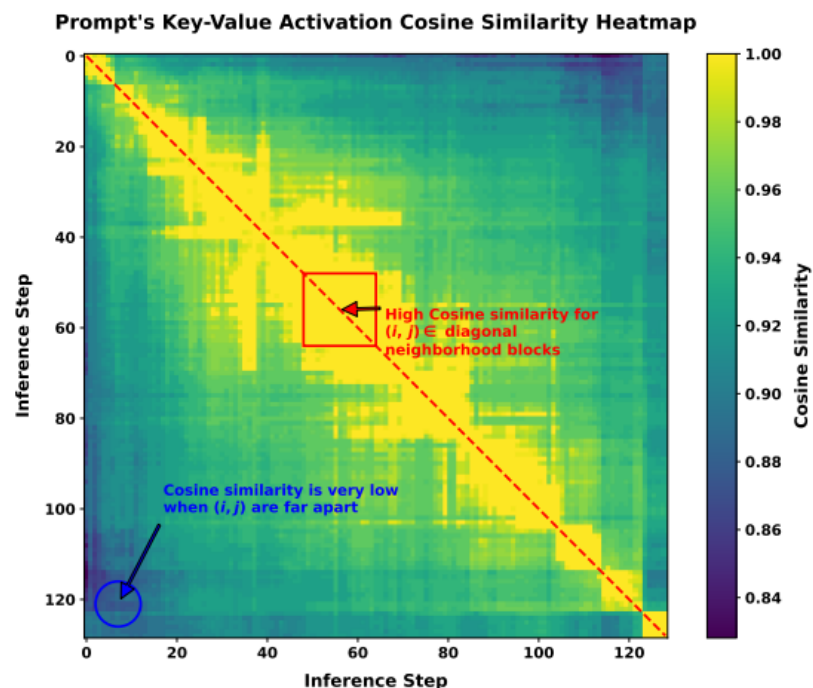


(b) **DualCache**: Bidirectional KV cache contains prefix and suffix Cache.

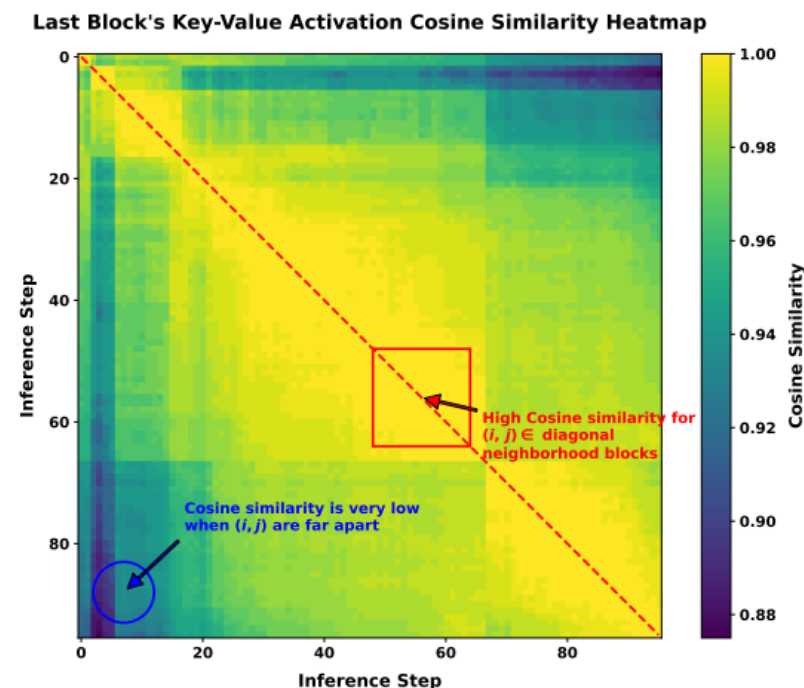
# DualCache Intuition

Why works?

- Observation: KV activations exhibit **high similarity** across **adjacent** inference steps **within a block**

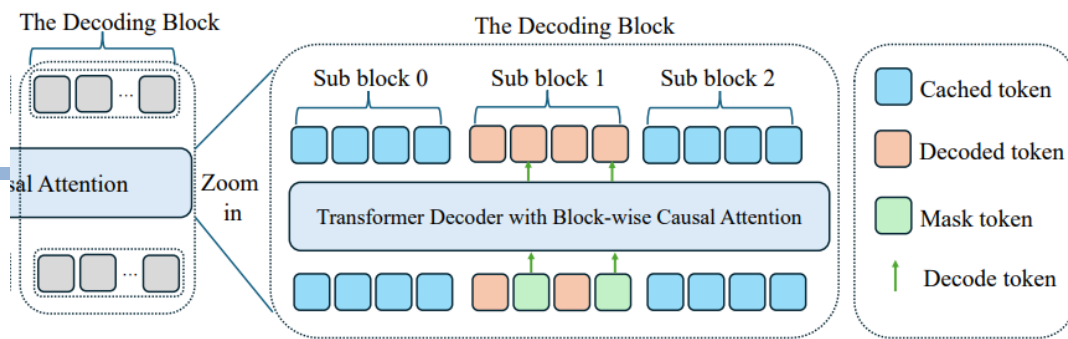


(a) Prompt block



(b) Last block

# Diffusion within blocks

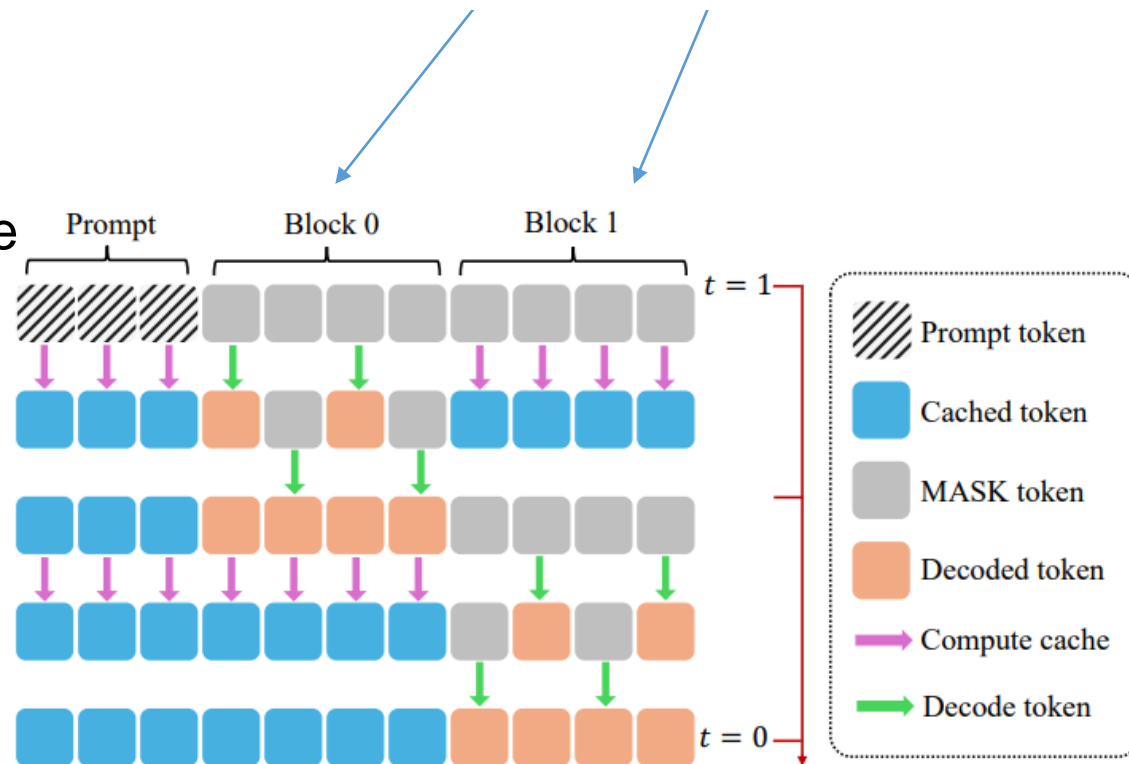


Cannot leverage KV cache?

❑ **DualCache**: approximate KV cache

❖ Cache previous kv in one sub block inference

❖ **Update** when current sub block is finished



(b) **DualCache**: Bidirectional KV cache contains prefix and suffix Cache.

# Performance under different cache block size

- ❑ Smaller block size incur **overhead**
- ❑ larger block size **diminish** accuracy
- ❑ Trade-off

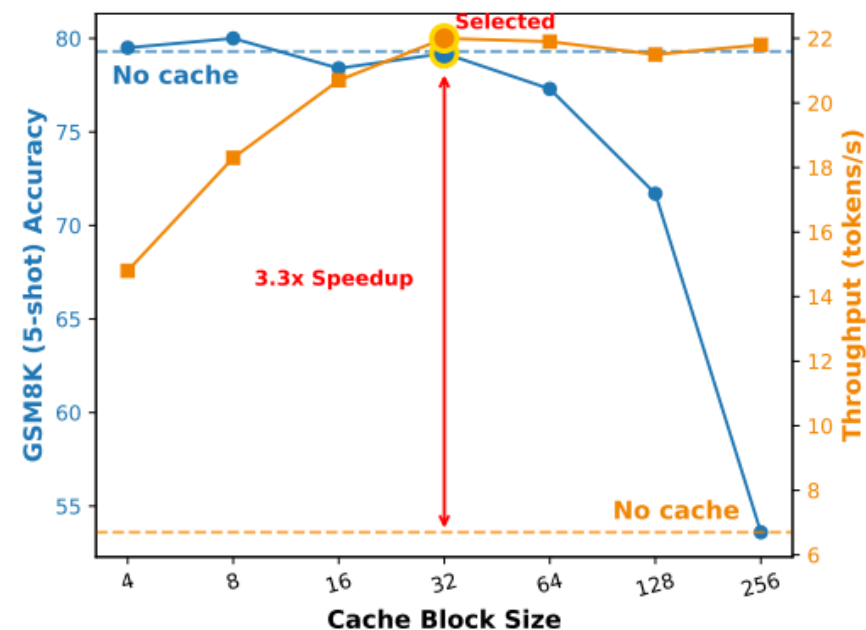
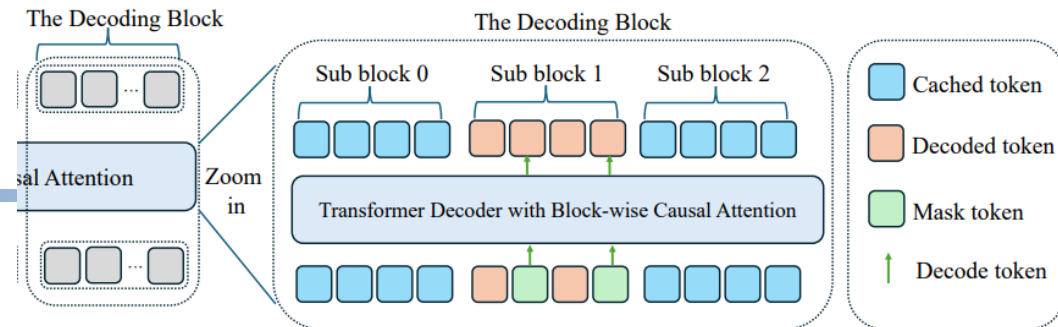


Figure 4 | **Impact of Cache Block Size on Accuracy and Throughput.** The orange line illustrates the effect of varying cache block size on throughput, while the blue line depicts accuracy.

# Diffusion within blocks



## □ Confidence-Aware Parallel Decoding

(TL;DR) If the model is confident on many positions, parallel decoding will not introduce errors.

**Theorem 1** (Parallel Decoding under High Confidence). *Suppose there exists a specific sequence of tokens  $\mathbf{x}^* = (x_{i_1}, \dots, x_{i_n})$  such that for each  $j \in \{1, \dots, n\}$ , the model has high confidence in  $x_{i_j}$ :  $p_j(X_{i_j} = x_{i_j} | E) > 1 - \epsilon$  for some small  $\epsilon > 0$ . Then, the following results hold:*

1. *Equivalence for Greedy Decoding: If  $(n + 1)\epsilon \leq 1$  (i.e.,  $\epsilon \leq \frac{1}{n+1}$ ), then*

$$\operatorname{argmax}_{\mathbf{z}} p(\mathbf{z} | E) = \operatorname{argmax}_{\mathbf{z}} q(\mathbf{z} | E) = \mathbf{x}^*. \quad (4)$$

*This means that greedy parallel decoding (selecting  $\operatorname{argmax} q$ ) yields the same result as greedy sequential decoding (selecting  $\operatorname{argmax} p$ ).*

*This bound is tight: if  $\epsilon > \frac{1}{n+1}$ , there exist distributions  $p(\mathbf{X} | E)$  satisfying the high-confidence marginal assumption for which  $\operatorname{argmax}_{\mathbf{z}} p(\mathbf{z} | E) \neq \operatorname{argmax}_{\mathbf{z}} q(\mathbf{z} | E)$ .*



# Diffusion within blocks

---

**Algorithm 1** Block-wise Confidence-aware Parallel Decoding with (Dual) KV Cache

---

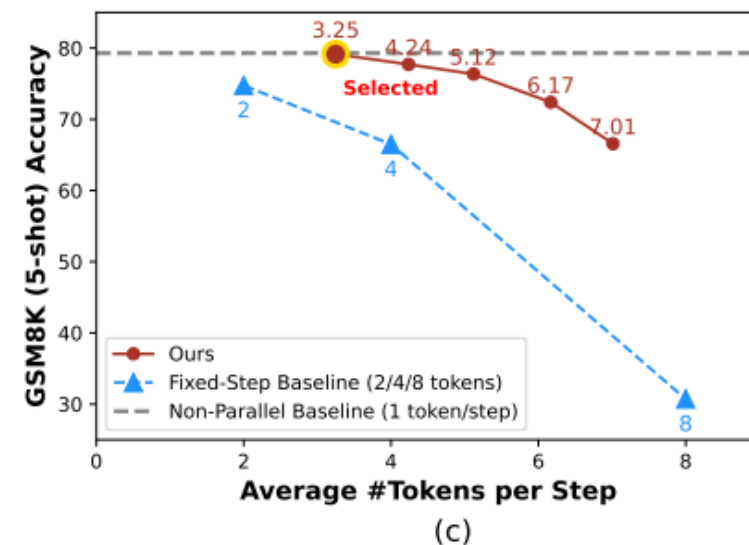
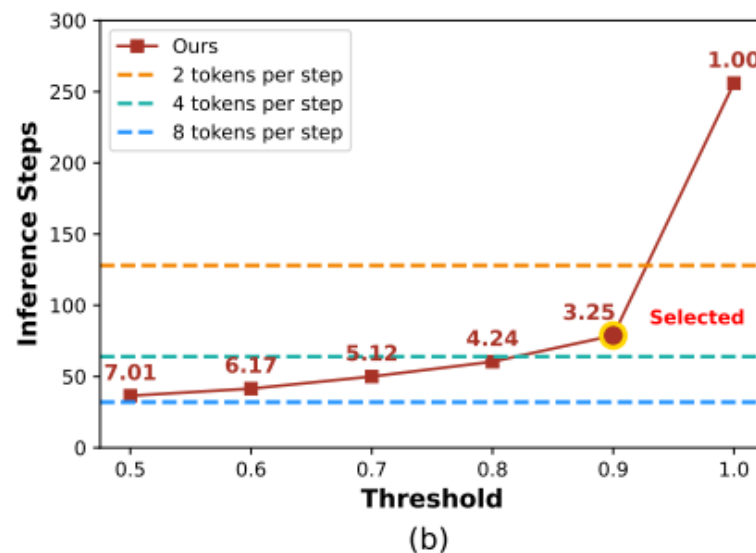
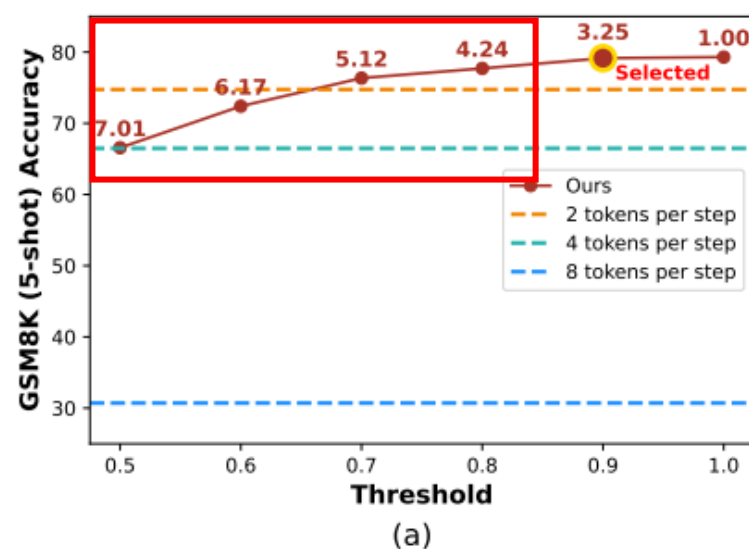
**Require:**  $p_\theta$ , prompt  $p_0$ , answer length  $L$ , blocks  $K$ , block size  $B$ , steps per block  $T$ , threshold  $\tau$ , use\_DualCache, strategy  $\in \{\text{threshold}, \text{factor}\}$ , factor  $f$

```
1:  $x \leftarrow [p_0; [\text{MASK}], \dots, [\text{MASK}]]$ 
2: Initialize KV Cache (single or dual) for  $x$  (fuse with decoding). // KV Cache Init
3: for  $k = 1$  to  $K$  do
4:    $s \leftarrow |p_0| + (k - 1)B$ ,  $e \leftarrow |p_0| + kB$ 
5:   for  $t = 1$  to  $T$  do
6:     Use cache, run  $p_\theta$  on  $x^{[s,e)}$  if use_DualCache else  $x^{[s,:)}$  // Cache Reuse
7:     For masked  $x^i$ , compute confidence  $c^i = \max_x p_\theta(x^i | \cdot)$  // Confidence scoring
8:     if strategy == threshold then
9:       Unmask all  $i$  in  $[s, e)$  with  $c^i \geq \tau$ , always unmask max  $c^i$ 
10:    else if strategy == factor then
11:      Sort  $c^i$  in descending order as  $(c^{(1)}, c^{(2)}, \dots, c^{(m)})$ 
12:      Find largest  $n$  such that  $(n + 1)(1 - c^{(n)}) < f$ 
13:      Unmask top- $n$  tokens, always unmask the max  $c^i$ 
14:    end if
15:    if all  $x^{[s,e)}$  unmasked then
16:      break
17:    end if
18:  end for
19:  Update KV cache: if use_DualCache: prefix & suffix; else: prefix. // Cache Update
20: end for
21: return  $x$ 
```

Unmask **dynamic** number of tokens

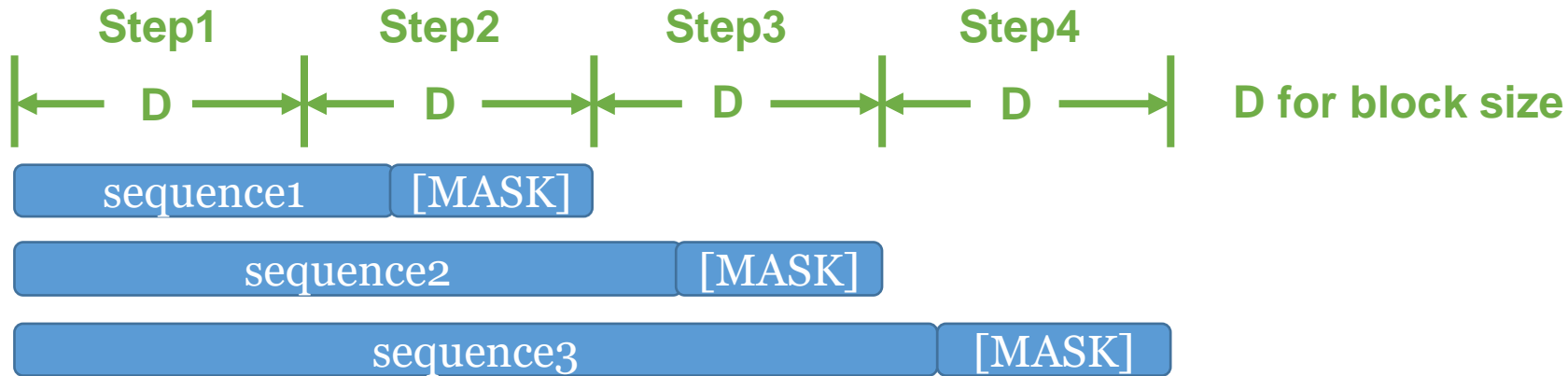


# Threshold Curve



- ❖ Baseline: llama previous fixed-N tokens
- ❖ Setting: GSM8K, 256 length

# Batch decoding with padding



❑ Adaptation to block diffusion

# Outline

---

- Background
- Design
- Evaluations**
- Discussion

# Setups

---

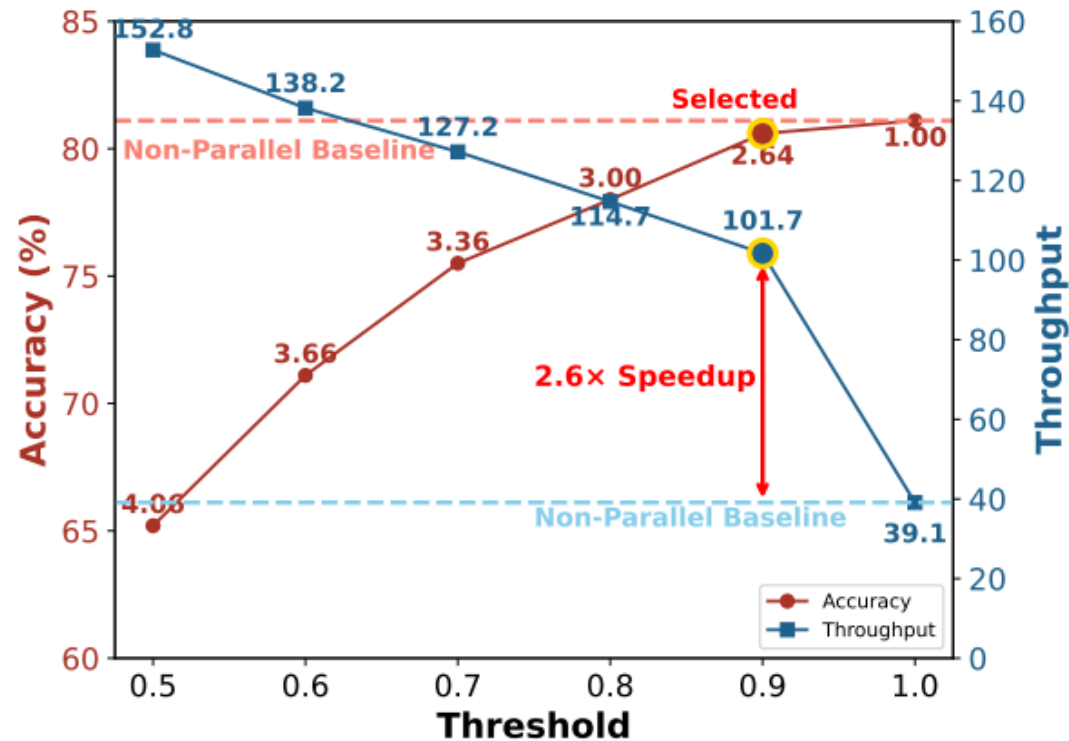
- ❑ Tuning model: Qwen-2.5 1.5B and 7B instruct
- ❑ Training dataset: LLaMA-Nemotron post-training dataset (batch 256)
- ❑ Training environment: 64 NVIDIA A100 GPUs
- ❑ Training configuration: 1.5B: learning rate  $2 \times 10^{-5}$  for 6,000 steps (costs 8h)  
7B: learning rate  $1 \times 10^{-5}$  for 2500 steps (costs 12h)
- ❑ Block size: 32
- ❑ Sub-Block size: 8

# Benchmark Results

| Model                | #Params | HumanEval |      | MBPP |      | GSM8K | Math | IFEval | MMLU | GPQA | Avg.        |
|----------------------|---------|-----------|------|------|------|-------|------|--------|------|------|-------------|
|                      |         | Base      | Plus | Base | Plus |       |      |        |      |      |             |
| 1B Models            |         |           |      |      |      |       |      |        |      |      |             |
| LlaMA-3.2            | 1.2B    | 34.1      | 31.1 | 34.1 | 29.4 | 43.0  | 23.8 | 58.9   | 44.4 | 24.1 | 35.9        |
| SmolLM 2             | 1.7B    | 34.1      | 28.7 | 50.6 | 46.0 | 47.7  | 21.1 | 55.1   | 49.1 | 29.2 | 40.7        |
| Qwen2.5-1.5B         | 1.5B    | 42.1      | 37.2 | 48.1 | 41.3 | 57.0  | 46.8 | 41.2   | 54.6 | 30.6 | <u>44.3</u> |
| Qwen2.5-1.5B-Nemo-FT | 1.5B    | 37.2      | 33.5 | 53.4 | 44.4 | 58.5  | 43.5 | 39.4   | 58.1 | 31.0 | <u>44.3</u> |
| Fast-dLLM v2         | 1.5B    | 43.9      | 40.2 | 50.0 | 41.3 | 62.0  | 38.1 | 47.0   | 55.1 | 27.7 | <b>45.0</b> |
| 7B+ Models           |         |           |      |      |      |       |      |        |      |      |             |
| LLaDA                | 8B      | 35.4      | 31.7 | 31.5 | 28.6 | 78.6  | 26.6 | 59.9   | 65.5 | 31.8 | 43.3        |
| LLaDA-1.5            | 8B      | 52.4      | -    | 42.8 | -    | 83.3  | 42.6 | 58.2   | 66.0 | 36.9 | -           |
| LLaDA-MoE            | 7B      | 61.6      | -    | 70.0 | -    | 82.4  | 58.7 | 59.3   | 67.2 | -    | -           |
| Dream                | 7B      | 57.9      | 53.7 | 68.3 | 56.1 | 81.0  | 39.2 | 62.5   | 67.0 | 33.0 | 57.6        |
| Qwen2.5-7B           | 7B      | 51.2      | 47.6 | 57.7 | 49.5 | 71.4  | 73.3 | 70.8   | 68.7 | 33.5 | 58.2        |
| Qwen2.5-7B-Nemo-FT   | 7B      | 52.4      | 48.2 | 57.1 | 50.0 | 84.1  | 72.0 | 69.5   | 68.6 | 34.2 | <u>59.6</u> |
| Fast-dLLM v2         | 7B      | 63.4      | 58.5 | 63.0 | 52.3 | 83.7  | 61.6 | 61.4   | 66.6 | 31.9 | <b>60.3</b> |

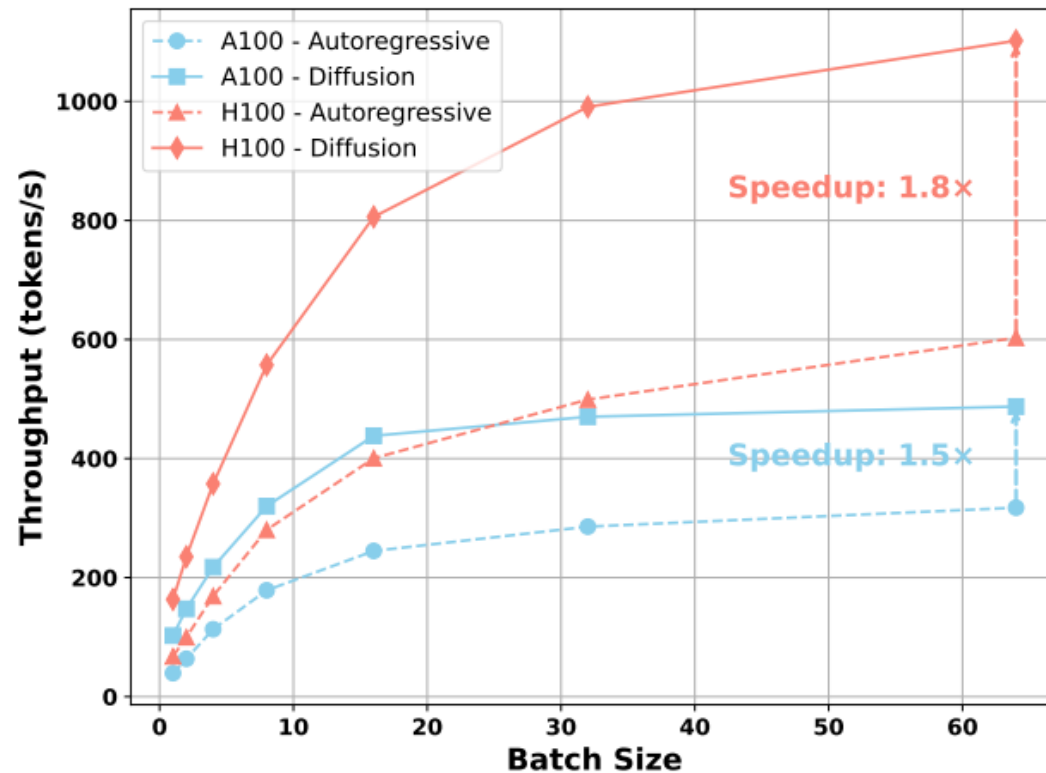
# Performance

□ Accuracy and throughput under different thresholds on GSM8K



# Performance

- ❑ **Throughput** comparison between AR(Qwen2.5-7B-Instruct) and diffusion(Fast-dLLM v2 7B) generation for GSM8K



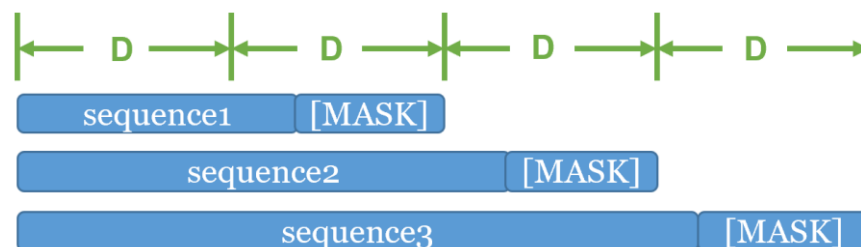
# Ablation Study

❖ Based on Fast-dLLM v2 1.5B

| Method            | HumanEval   |             | MBPP        |             | GSM8K       | Math        | IFEval      | MMLU        | GPQA        | Avg.        |
|-------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|                   | Base        | Plus        | Base        | Plus        |             |             |             |             |             |             |
| Naive token shift | <u>38.4</u> | 32.9        | 44.4        | <u>38.6</u> | 59.0        | <u>37.3</u> | 39.9        | 52.9        | <b>27.9</b> | 41.3        |
| + pad             | <u>38.4</u> | <u>34.1</u> | <u>45.2</u> | 38.4        | <u>60.1</u> | 37.0        | <u>45.8</u> | <u>53.5</u> | <u>27.7</u> | <u>42.2</u> |
| + pad + CM        | <b>43.9</b> | <b>40.2</b> | <b>50.0</b> | <b>41.3</b> | <b>62.0</b> | <b>38.1</b> | <b>47.0</b> | <b>55.1</b> | <u>27.7</u> | <b>45.0</b> |

❖ Naive token shift: use hidden state at  $i-1$  to predict  $x_i$

❖ Pad:



❖ CM: complementary mask



# Ablation Study

## ❑ Sub-Block size and Block size affect performance

Table 3 | Sub-Block size decoding improves performance, with size 8 being optimal.

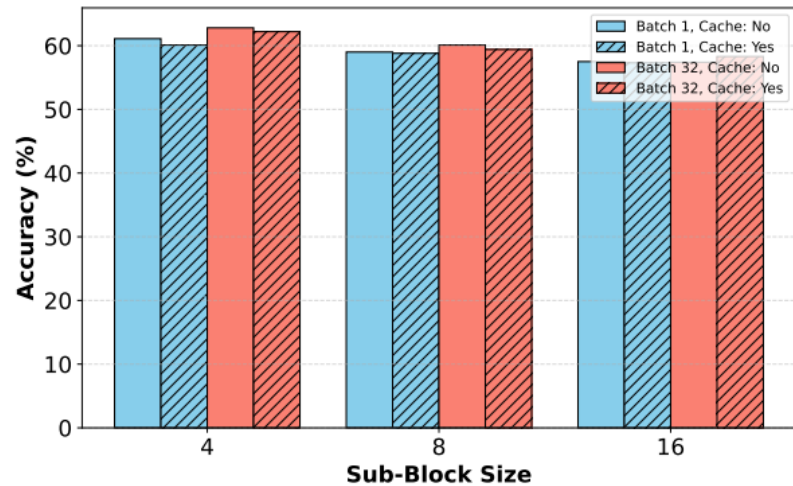
| Sub-Block Size | 2           | 4           | 8           | 16   | 32   |
|----------------|-------------|-------------|-------------|------|------|
| GSM8K          | <b>62.8</b> | 61.8        | <u>62.0</u> | 61.3 | 60.2 |
| HumanEval      | 42.7        | <u>43.3</u> | <b>43.9</b> | 39.6 | 38.4 |
| HumanEval+     | <u>39.6</u> | <b>40.2</b> | <b>40.2</b> | 36.0 | 34.8 |

Table 4 | Inference with mismatched sizes reduces performance.

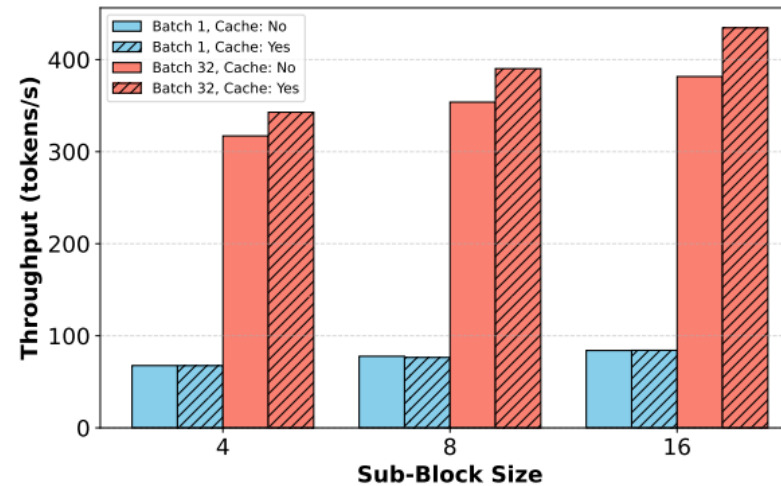
| Block Size | 2    | 4           | 8           | 16          | 32          |
|------------|------|-------------|-------------|-------------|-------------|
| GSM8K      | 53.2 | 56.8        | 58.5        | <u>59.7</u> | <b>60.2</b> |
| HumanEval  | 37.8 | <b>43.3</b> | <b>43.3</b> | <u>38.4</u> | <u>38.4</u> |
| HumanEval+ | 34.1 | <u>39.0</u> | <b>39.6</b> | 34.1        | 34.8        |

# Ablation Study

❑ Cache is a purely efficiency-enhancing feature



(a)



(b)

Figure 6 | Effect of small block size and sub-block cache on model performance. **(a)** Accuracy remains largely unaffected by the use of sub-block cache across different block sizes and batch sizes. **(b)** Throughput increases as small block size grows due to higher decoding parallelism. While sub-block cache has negligible effect when batch size is small, it significantly improves throughput under compute-bound settings (e.g., batch size = 32).

# Outline

---

- Background
- Design
- Evaluations
- Discussion**

# Discussion

---

## ❑ Trade-off between accuracy and throughput in MDMs

- ❖ Parallel generation inevitably introduces conditional independence

## ❑ Arbitrary-Order Autoregressive?

- ❖ Semi-AR Semi-diffusion

## ❑ Data-efficient fine-tuning

- ❖ Fast-dLLM v2 achieves lossless adaptation with just  $\sim 1$ B tokens, compared to  $\sim 500$ B tokens required by Dream
- ❖ Diffusion Beats Autoregressive in Data-Constrained Settings (NeurIPS 25)

**Thanks!**