

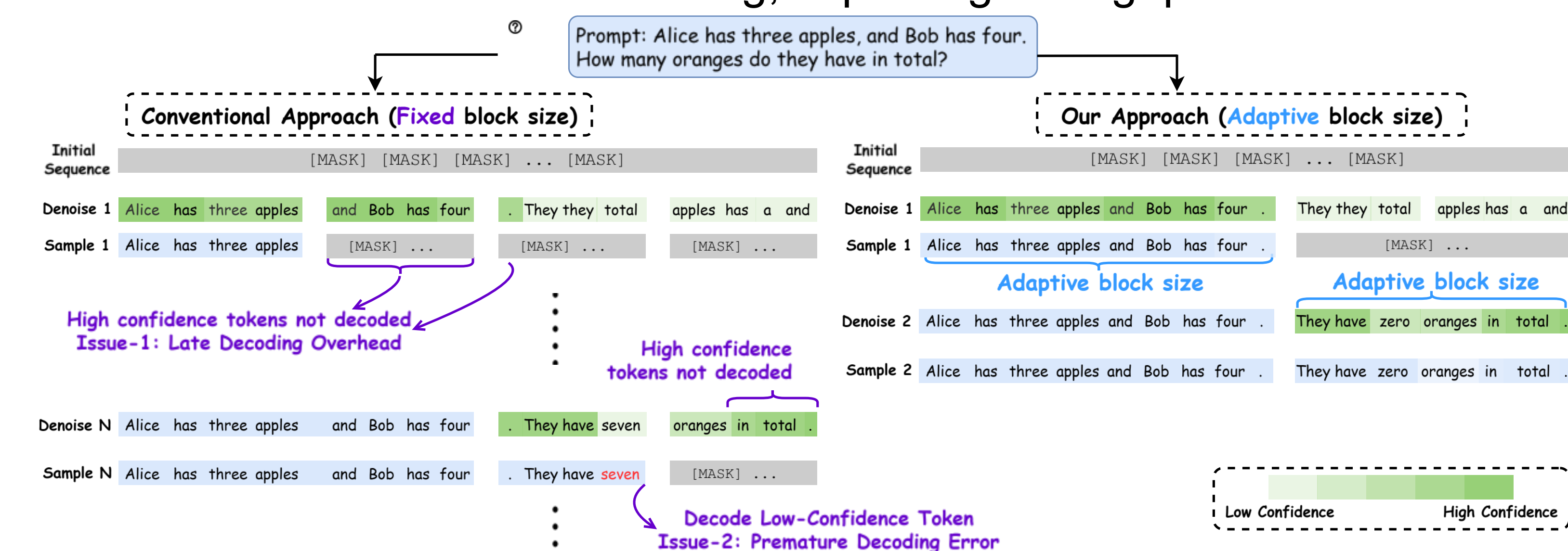
## Background & Motivation

### ➤ Diffusion-Based Large Language Models (dLLMs)

- Iteratively denoise a sequence with <MASK> tokens.
- Emerge as a promising alternative to AR models with parallel decoding and better controllability.

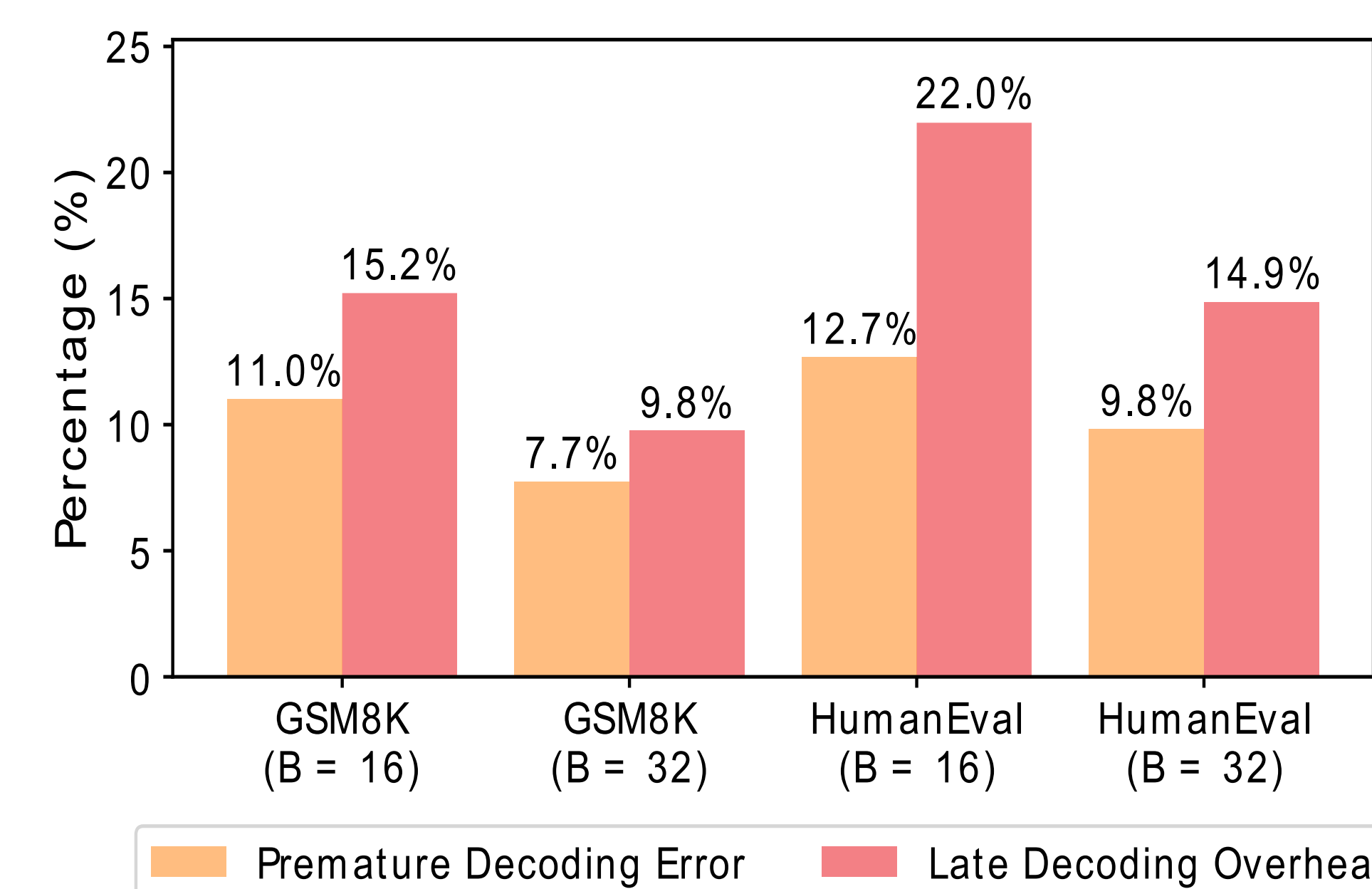
### ➤ Confidence-Based Semi-AR Decoding

- dLLM inference throughput is typically constrained by the lack of KV caching and the parallelism curse.
- Confidence-based semi-AR decoding balances parallelism and quality, and allows for block-level caching, improving throughput.



### ➤ Fundamental Issues of Fixed Block Size

- Late Decoding Overhead:** Delays the unmasking of high-confidence tokens outside the current block (throughput ↓).
- Premature Decoding Error:** Forces early commitment to low-confidence tokens within each block (quality ↓).



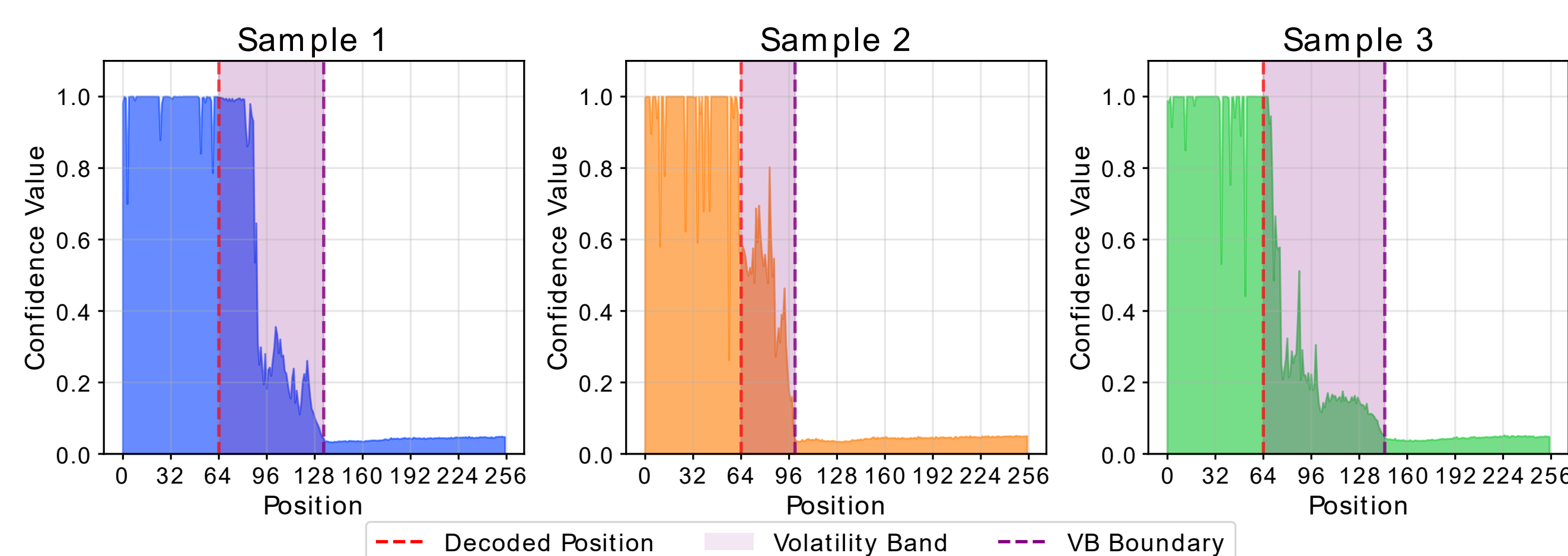
## Method

### ➤ Confidence Dynamics

- A high-confidence region emerges near the decoded token.
- The decoding trace displays a global autoregressive sampling tendency driven by semantic dependencies.
- The high-confidence region extends toward adjacent positions; other positions maintain low confidence.

### ➤ Volatility Band and Local Stochasticity

- Partitioning the confidence landscape: high-confidence plateau, volatility band (VB), and low-confidence floor.
- Confidence pattern in VB is locally stochastic and dependent on the semantic context.



### ➤ Semantic-Aware Block Size Scheduling

- A scheduling mechanism that sets the block size guided by the length of the current semantic step
- Use confidence of a set of delimiter tokens to represent the semantic boundary and predict the block size.

#### Algorithm 1 Semantic-Aware Block Size Determination

**Inputs:** predicted sequence  $\hat{y}$ ; confidences  $c$ ; generation budget  $L$ ; default block size  $B_0$ ; delimiter set  $\mathcal{D}$ ; delimiter threshold  $\tau_D$ ; current position  $g$ .

**Output:** block size  $B$

```

1: function COMPUTEBLOCKLENGTH( $\hat{y}$ ,  $c$ ,  $L$ ,  $B_0$ ,  $\mathcal{D}$ ,  $\tau_D$ ,  $g$ )
2:    $\triangleright$  Sampling window boundary
3:    $start, remaining \leftarrow g, L - g$ 
4:    $w \leftarrow \min(\max(1, \lfloor 0.25 \cdot g \rfloor), remaining)$ 
5:    $W \leftarrow \{start, \dots, start + w - 1\}$   $\triangleright$  window token indices
6:    $\triangleright$  Find highest-confidence delimiter
7:    $\mathcal{I} \leftarrow \{i \in W \mid \hat{y}_i \in \mathcal{D}\}$ 
8:   if  $\mathcal{I} \neq \emptyset$  then
9:      $pos \leftarrow \arg \max_{i \in \mathcal{I}} c_i$   $\triangleright$  Select position with max delimiter token confidence
10:     $c_{max} \leftarrow c_{pos}$ 
11:  else
12:     $c_{max} \leftarrow -\infty$ 
13:  end if
14:   $\triangleright$  Determine block size
15:  if  $c_{max} \geq \tau_D$  then
16:     $B \leftarrow (pos - start + 1)$   $\triangleright$  inclusive length up to the delimiter token
17:  else
18:     $B \leftarrow \min(B_0, remaining)$ 
19:  end if
20:  return  $B$ 
21: end function

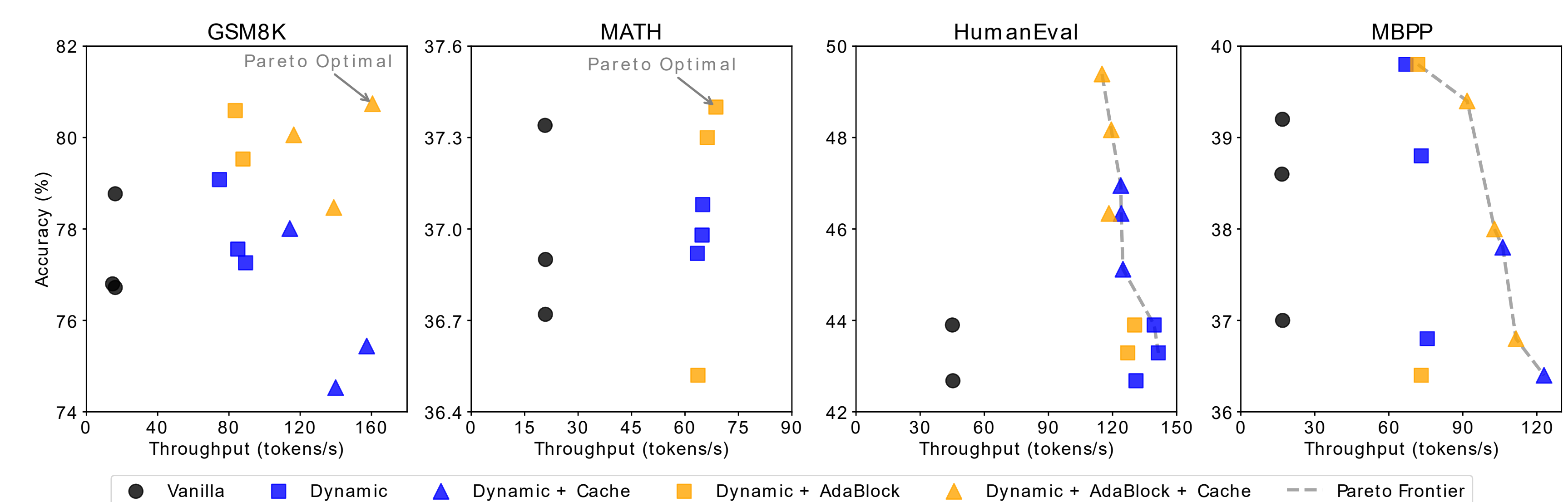
```

## Experiments

### ➤ Main Results

Method	LLaDA-Instruct			LLaDA-1.5			Dream-Base		
	$B_0 = 16$	$B_0 = 32$	$B_0 = 64$	$B_0 = 16$	$B_0 = 32$	$B_0 = 64$	$B_0 = 16$	$B_0 = 32$	$B_0 = 64$
<b>GSM8K</b>									
Vanilla	78.8	76.7	76.8	82.3	82.3	80.4	76.3	76.4	75.1
Dynamic	79.1	77.6	77.3	82.6	82.2	80.7	75.5	75.5	75.6
+Ada	80.6 <sup>+1.5</sup>	80.6 <sup>+3.0</sup>	79.5 <sup>+2.2</sup>	83.0 <sup>+0.4</sup>	82.4 <sup>+0.2</sup>	80.3 <sup>-0.4</sup>	75.7 <sup>+0.2</sup>	75.7 <sup>+0.2</sup>	75.9 <sup>+0.3</sup>
+Cache	78.0	74.5	75.4	80.7	80.2	80.0	75.6	74.5	74.6
+Ada+Cache	80.0 <sup>+2.0</sup>	78.5 <sup>+4.0</sup>	80.7 <sup>+5.3</sup>	81.3 <sup>+0.6</sup>	81.7 <sup>+1.5</sup>	79.7 <sup>-0.3</sup>	76.5 <sup>+0.9</sup>	75.1 <sup>+0.6</sup>	74.6 <sup>+0.0</sup>
<b>HumanEval</b>									
Vanilla	43.9	43.9	42.7	39.0	36.6	38.4	53.7	52.4	54.3
Dynamic	42.7	43.9	43.3	36.6	37.8	36.6	53.0	51.2	52.4
+Ada	43.3 <sup>+0.6</sup>	43.3 <sup>-0.6</sup>	43.9 <sup>+0.6</sup>	37.8 <sup>+1.2</sup>	38.4 <sup>+0.6</sup>	38.4 <sup>+1.8</sup>	53.7 <sup>+0.7</sup>	51.2 <sup>+0.0</sup>	53.7 <sup>+1.3</sup>
+Cache	45.1	46.3	47.0	33.5	36.0	34.1	50.0	53.0	56.1
+Ada+Cache	49.4 <sup>+4.3</sup>	46.3 <sup>+0.0</sup>	48.2 <sup>+1.2</sup>	36.0 <sup>+2.5</sup>	39.0 <sup>+3.0</sup>	36.0 <sup>+1.9</sup>	52.4 <sup>+2.4</sup>	53.0 <sup>+0.0</sup>	57.3 <sup>+1.2</sup>
<b>MATH</b>									
Vanilla	36.7	36.9	37.3	36.3	37.0	34.4	39.8	40.2	40.1
Dynamic	37.0	36.9	37.1	36.3	36.7	34.4	39.7	39.9	39.9
+Ada	36.5 <sup>-0.5</sup>	37.3 <sup>+0.4</sup>	37.4 <sup>+0.3</sup>	36.8 <sup>+0.5</sup>	36.7 <sup>+0.0</sup>	34.1 <sup>-0.3</sup>	39.6 <sup>-0.1</sup>	39.9 <sup>+0.0</sup>	39.9 <sup>+0.0</sup>
+Cache	35.4	35.8	36.0	34.9	33.2	32.1	38.0	38.5	38.8
+Ada+Cache	35.8 <sup>+0.4</sup>	35.3 <sup>-0.5</sup>	35.6 <sup>-0.4</sup>	35.2 <sup>+0.3</sup>	33.9 <sup>+0.7</sup>	32.4 <sup>+0.3</sup>	37.8 <sup>-0.2</sup>	38.4 <sup>-0.1</sup>	38.4 <sup>-0.4</sup>
<b>MBPP</b>									
Vanilla	39.2	38.6	37.0	38.2	37.0	23.2	12.4	12.4	12.8
Dynamic	39.8	38.8	36.8	38.2	37.0	24.6	12.6	12.4	12.2
+Ada	40.2 <sup>+0.4</sup>	39.8 <sup>+1.0</sup>	36.4 <sup>-0.4</sup>	39.4 <sup>+1.2</sup>	37.6 <sup>+0.6</sup>	29.8 <sup>+5.2</sup>	12.8 <sup>+0.2</sup>	14.2 <sup>+1.8</sup>	12.4 <sup>+0.2</sup>
+Cache	35.6	37.8	36.4	38.0	34.8	19.8	12.8	11.6	9.6
+Ada+Cache	39.4 <sup>+3.8</sup>	38.0 <sup>+0.2</sup>	36.8 <sup>+0.4</sup>	36.6 <sup>-1.4</sup>	36.4 <sup>+1.6</sup>	36.6 <sup>+6.8</sup>	12.8 <sup>+0.0</sup>	11.6 <sup>+0.0</sup>	12.4 <sup>+2.8</sup>

### ➤ Accuracy (%) VS Throughput (TPS)



### ➤ Ablation Studies

Model	$\tau_D = 0.3$	$\tau_D = 0.5$	$\tau_D = 0.7$
LLaDA-Instruct	<b>80.59</b>	79.08	77.94
Dream-Base	75.66	<b>75.74</b>	<b>75.74</b>

Ablation on Delimiter Threshold			
Method	$B_0 = 16$	$B_0 = 32$	$B_0 = 64$
Vanilla	69.1	66.7	61.3
Dynamic	69.0	66.7	61.2
+Ada	68.4 <sup>-0.6</sup>	67.5 <sup>+0.8</sup>	64.4 <sup>+3.2</sup>
+Cache	67.5	64.6	59.4
+Ada+Cache	68.9 <sup>+1.4</sup>	66.4 <sup>+1.8</sup>	62.7 <sup>+3.3</sup>

Delimiter Set	Acc. (%)
None (+Cache)	74.5
$\{\{\backslash n\}\}$	<b>78.5</b>
$\{\{, \}\}$	75.1
$\{\{. \}\}$	74.5
$\{\{, \}, [ \cdot ]\}$	75.1
$\{\{\backslash n\}, [ \cdot ]\}$	<b>78.5</b>
$\{\{\backslash n\}, [ \cdot ]\}$	<b>78.3</b>
$\{\{\backslash n\}, [ \cdot ], [ \cdot ]\}$	<b>78.7</b>

Accuracy on IFEval

Ablation on Delimiter Set