

# Plan for Speed

Dilated Scheduling for Masked Diffusion Language Models

Omer Luxembourg ▪ Haim Permuter ▪ Eliya Nachmani

*School of Electrical and Computer Engineering, Ben-Gurion University of the Negev*



**ICML**  
International Conference  
On Machine Learning



# Motivation

- Masked diffusion LMs promise **parallel generation**.
- But existing confidence and entropy based planners, in high probability, **pick neighbouring tokens**.
- Due to the model architecture, they ignore dependencies between co-decoded positions.
- Effectively yielding **near-autoregressive, clustered decoding**.

*Inference planners can reach different quality and speedups; DUS is an inference planner.*

# Background: masked diffusion language models

- **Forward process:** each token is replaced by [MASK] with probability  $t$ .
- **Training:** predict masked tokens with reweighted cross-entropy

$$\mathcal{L}(\theta) = -\mathbb{E}_{t, \mathcal{X}, \mathcal{M}_t} \left[ \frac{1}{t} \sum_{i: M_{t,i} = [\text{MASK}]} \log p_{\theta}(X_i | \mathcal{S}_t) \right].$$

- **Inference loop:** *planner* selects masked positions  $\rightarrow$  denoiser predicts them  $\rightarrow$  state updates.
- **Key lever:** planners control the quality and speed of generation. **DUS is a planner.**

# Dilated Unmasking Scheduler (DUS)

Within each block of  $B$  tokens, at iteration  $t = 1, \dots, R$ , DUS unmasking positions

$$\mathcal{I}_t = \left\{ k \in \{1, \dots, B\} \setminus \mathcal{U}_{t-1} \mid (k-1) \bmod s_t = 0 \right\},$$

with

$$s_t = \lfloor B/2^t \rfloor, \quad R = \lceil \log_2 B \rceil, \quad \mathcal{U}_t = \mathcal{U}_{t-1} \cup \mathcal{I}_t, \quad \mathcal{U}_0 = \emptyset.$$

1/8	DUS	[M]	[M]	[M]	[M]	[M]	[M]	[M]	we	[M]	[M]	[M]	[M]	[M]	[M]	[M]
2/4	DUS	[M]	[M]	[M]	are	[M]	[M]	[M]	we	[M]	[M]	[M]	in	[M]	[M]	[M]
3/2	DUS	[M]	tokens	[M]	are	[M]	apart	[M]	we	[M]	fill	[M]	in	[M]	log-B	[M]
4/1	DUS	picks	tokens	that	are	far	apart	so	we	can	fill	them	in	just	log-B	steps

$B=16 \Rightarrow 4$  denoiser calls ( $s_t=8, 4, 2, 1$ ). **dark**=newly unmasked, **light**=already unmasked, white=masked.

# Novelties

- **Inference-only, model-agnostic.** Drop-in scheduler; zero changes to architecture or training.
- **Logarithmic schedule.** Deterministic coarse-to-fine dilation in  $\mathcal{O}(\log B)$  denoiser calls per block.
- **Theory.** The cheap sum-of-marginals objective is tightest exactly when spacing is largest.
- **Empirical breadth.** 8 benchmarks  $\times$  5 models (LLaDA, Dream, DiffuCoder).
- **Hybrid extension.** Dilated spacing as a post-filter on top of adaptive EB/CB samplers.

## Theory: VLMC spacing bound (informal)

Assume a stationary, ergodic **fast-mixing** chain. For positions in the same DUS step with spacing  $s_t$ :

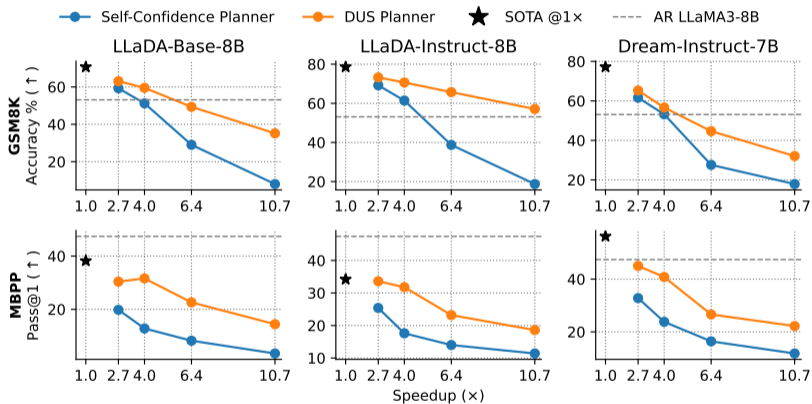
$$I(X_i; X_{i+s_t}, X_{i+2s_t}, \dots) \leq C \rho^{s_t}, \quad \rho \in (0, 1).$$

Mutual information decays *exponentially* with spacing, so the gap between the true joint entropy and the cheap sum-of-marginals shrinks:

$$H(X_{\mathcal{I}_t} | \mathcal{S}_t) \geq \sum_{i \in \mathcal{I}_t} H(X_i | \mathcal{S}_t) - \varepsilon(s_t), \quad \varepsilon(s_t) \text{ exp. small in } s_t.$$

*Tightest at coarse spacing; later steps lean on the richer revealed context.*

# Main result: the speed-quality frontier



DUS (orange) vs self-confidence (blue) vs Llama-3-8B AR (gray dashed), on GSM8K, MATH500, HumanEval, MBPP. Up to +27 points over self-confidence at matched compute.

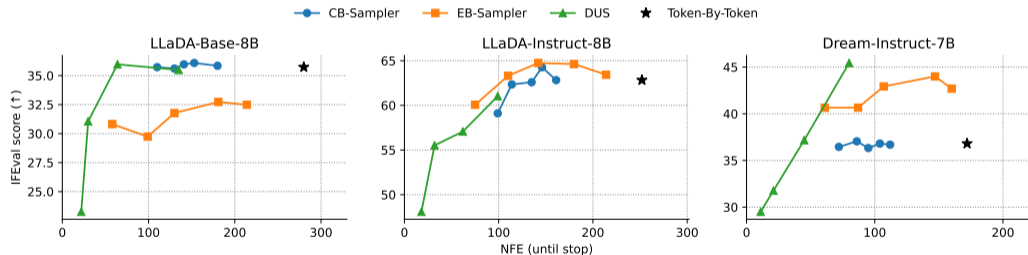
# Consistent gains across models and tasks

Data	Model	TBT	$B=16 (\times 4)$		$B=32 (\times 6.4)$		$B=64 (\times 10.7)$	
		Conf.	Conf.	DUS	Conf.	DUS	Conf.	DUS
<b>GSM8K</b>	LLaDA-Base	72.6	51.2	<b>59.5</b>	29.0	<b>49.4</b>	8.0	<b>35.2</b>
	LLaDA-Inst	80.3	61.4	<b>70.7</b>	38.7	<b>65.7</b>	18.7	<b>57.1</b>
	Dream-Inst	77.1	53.2	<b>56.6</b>	27.6	<b>44.7</b>	17.9	<b>32.1</b>
<b>MBPP</b>	LLaDA-Base	38.0	12.8	<b>31.6</b>	8.2	<b>22.6</b>	3.4	<b>14.4</b>
	LLaDA-Inst	39.4	17.6	<b>31.8</b>	14.0	<b>23.2</b>	11.4	<b>18.6</b>
	Dream-Inst	56.2	23.8	<b>40.8</b>	16.4	<b>26.6</b>	11.8	<b>22.2</b>

GSM8K accuracy / MBPP pass@1 (%). **Bold** = better planner. TBT = token-by-token baseline. The gap widens with speed-up: at  $\times 10.7$ , DUS more than **quadruples** self-confidence on GSM8K and **quadruples** it on MBPP (LLaDA-Base).

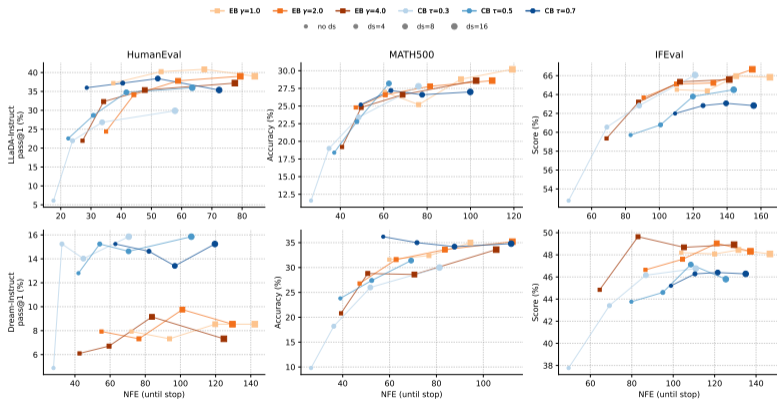
# Versus adaptive samplers (EB / CB)

On IFEval, DUS (**green**) sits **left** of entropy-bounded (EB) and confidence-bounded (CB) at every accuracy level: it reaches near token-by-token accuracy at 3–4× **fewer** denoiser calls (NFE), and on math/code DUS at  $B=16$  stays within  $\sim 2$  points of EB at fewer calls.



# Hybrid: dilated spacing as a post-filter

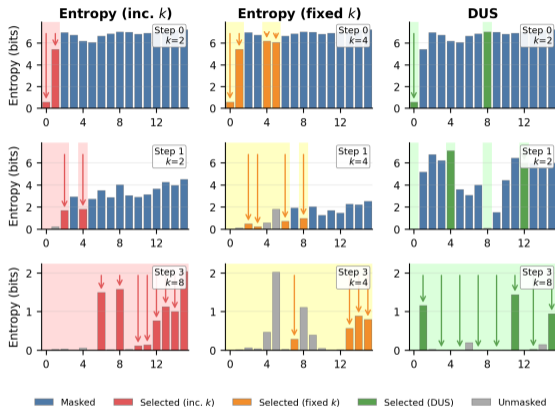
Apply dilated spacing *on top of* EB ( $\gamma=2$ ) and CB ( $\tau=0.5$ ) samplers, with no change to the score function: keep their scores, but enforce a minimum gap between co-decoded positions. Accuracy rises on **every** dataset and sampler for a modest NFE increase; biggest jump +13.4 pts on HumanEval (EB).



LLaDA-Inst (top) and Dream-Inst (bottom). Accuracy vs. NFE; marker size = spacing block (ds=4, 8, 16).

+spacing pushes every curve up-and-right.

# Why it works: DUS spaces, others cluster



Per-position entropy and each planner's picks at  $B=16$ .  
DUS (right) keeps co-decoded tokens  $\sim 2-3\times$  further apart.

The same spacing shows up in the statistics (LLaDA-Base):

Data	$B$	Avg. pairwise distance			
		DUS	SC	EB	CB
MATH500	16	<b>6.2</b>	3.0	3.5	4.5
MATH500	32	<b>9.6</b>	4.1	4.7	7.1
HumanEval	16	<b>6.2</b>	2.8	4.4	4.7
HumanEval	32	<b>11.6</b>	3.7	6.4	8.6

- DUS spreads co-decoded tokens  $2-3\times$  wider than SC / EB / CB.
- 100% of DUS picks are *isolated* (no adjacent co-decoded token) vs  $\sim 50\%$  for SC.
- Exactly the near-independent regime our theory calls safe.

# Takeaways

- *Space, don't cluster*: separate co-decoded tokens.
- Training-free, deterministic,  $\mathcal{O}(\log B)$  planner.
- Up to **+27** pts at matched compute; 3–4× fewer NFE than adaptive samplers.
- Composes as a post-filter on top of EB / CB.



[omerlux.github.io/DUS](https://omerlux.github.io/DUS)

Code & paper