

Fast-Decoding Diffusion Language Models via Progress-Aware Confidence Schedules

Amr Mohamed^{1,2*}, Yang Zhang², Michalis Vazirgiannis^{1,2}, Guokan Shang¹

¹MBZUAI, ²Ecole Polytechnique

Abstract

Diffusion large language models (dLLMs) offer a promising alternative to autoregressive models, but their practical utility is severely hampered by slow, iterative sampling. We present *SchED*, a training-free, model-agnostic early-exit algorithm that terminates diffusion decoding using a progress-aware confidence threshold. We evaluate *SchED* across multiple diffusion model families and a diverse set of benchmarks spanning multiple-choice, math, long-form QA, and translation. *SchED* delivers substantial acceleration: on instruction-tuned models, it achieves approximately $4\times$ speedups while retaining baseline performance on average. On base models, *SchED* yields consistent speedup gains with 99.1–100% performance retention, with up to $2.34\times$ under more aggressive settings. Under a conservative quality-penalized speed metric, *SchED* consistently outperforms prior confidence-based early-exit methods, including on long-form generation where existing approaches tend to break down. An entropy analysis of the model’s token predictions reveals that instruction tuning speeds up the decay of predictive entropy. By leveraging inherent confidence stabilization as a signal for computational efficiency, *SchED* provides a robust framework for efficient dLLM inference.

1 Introduction

Large Language Models (LLMs) have evolved rapidly in recent years, but the dominant decoding paradigm remains autoregressive (AR), a process that is inherently sequential and constrains opportunities for parallel generation and global context use (Brown et al., 2020; Yin et al., 2024; Zhang et al., 2025b). In response, Diffusion Large Language Models (dLLMs) have emerged as a credible alternative to AR decoding, offering compelling advantages such as parallel refinement, flexible

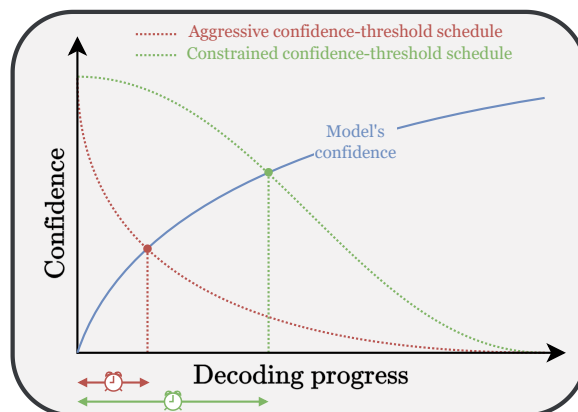


Figure 1: Progress-aware early exit in *SchED*. Decoding terminates when the model’s confidence (blue) intersects a progress-dependent threshold. The aggressive schedule (red) relaxes the threshold rapidly to trigger an earlier exit and maximize speedup. Conversely, the constrained schedule (green) maintains a higher threshold for a longer duration, ensuring greater output stability by delaying the exit point.

infilling, and bidirectional attention over the partially generated sequence (Zou et al., 2023). This paradigm has matured rapidly, with efforts to train powerful base and instruct models from scratch (Ye et al., 2025), adapt existing AR checkpoints (Gong et al., 2024b; Nie et al., 2025), and push capabilities into complex domains like reasoning and planning, where they can outperform AR models on certain tasks (Phuong et al., 2025; Gong et al., 2024a; Song et al., 2025; Pal et al., 2025), and principled scaling analyses for masked diffusion objectives (Li et al., 2024). Collectively, these results position dLLMs as a practical family of foundation-model architectures with distinct decoding affordances.

Despite this promise, decoding efficiency remains a central bottleneck for dLLMs. dLLM generation proceeds via a reverse-diffusion chain with many refinement steps. In addition, practitioners must choose step budgets and transfer schedules a

*Correspondence: amr.mohamed@mbzuai.ac.ae

priori, often conservatively, to avoid quality loss. This leads to unnecessary computation on “easy” inputs and, when heuristic decoding parameter choices (e.g., step budgets) are too extreme, results in unstable behavior across tasks. A growing body of work addresses this bottleneck through training-free early-commit methods, which exploit the empirical observation that predictions often stabilize well before the final diffusion step (Pengxiang et al., 2025). Dynamic policies modulate exploration versus acceleration across the chain, using signals like historical logits or local determinism to reduce redundant steps (Wei et al., 2025; Wang et al., 2025b; Or-El et al., 2025). Other lines of work focus on error correction and refinement, enabling models to revise their own outputs by re-masking low-confidence tokens (Zhang et al., 2025a; Su et al., 2025). Few- and one-step routes also address this bottleneck by transferring ideas from continuous diffusion via curriculum, consistency distillation, or flow matching (Sahoo et al., 2025; Chen et al., 2025; Meister et al., 2024). Orthogonal efforts reduce per-step latency with caching adapted to bidirectional refinement (Yuan et al., 2025b; Wang et al., 2025a) and with speculative mechanisms that draft and verify tokens in parallel (Sohoni et al., 2025; Yuan et al., 2025a). While impactful, these approaches often require additional training, introduce auxiliary models, or rely on complex heuristics, leaving room for a simple, training-free, and *architecture-agnostic* early-exit principle.

We frame diffusion decoding as a *when-to-stop* problem and introduce *SchED*, a training-free, model-agnostic early-exit mechanism. *SchED* aggregates token-level confidence (top-2 logit margins) over the answer region and compares it to a smooth, non-increasing threshold defined over normalized diffusion progress. This progress-aware criterion decouples stopping from fixed step budgets and terminates decoding once predictions have stabilized, avoiding the brittleness of discrete phase changes. *SchED* composes with standard transfer schedules (single-suffix or block diffusion) and requires no changes to model training (Nie et al., 2025; Ye et al., 2025). Figure 1 illustrates *SchED*’s progress-aware early-exit mechanism, showing how aggregated confidence interacts with different threshold schedules to trigger an exit.

We evaluate *SchED* across two diffusion LLM families (a single-block Dream decoder and a block-diffusion LLaDA decoder), each in base and instruction-tuned variants, on ten diverse bench-

marks covering multiple-choice, math, long-form QA/summarization, and translation. On instruction-tuned models, *SchED* retains 99.8–100% of baseline quality on average, while yielding $3.8\text{--}4.0\times$ speedups and outperforming recent training-free early-exit methods on a conservative quality-penalized speed metric; on base models, it delivers smaller but consistent gains at near-parity quality (Section 5). Together, these results show that *SchED* can substantially reduce diffusion decoding costs without sacrificing quality, and that instruction tuning lowers the confidence barrier for early exit in QA-style tasks, making it particularly effective in that regime. This paper makes the following contributions:

1. **Training-free, schedule-based early exit.** We introduce *SchED* (**S**chedule-based **E**arly **D**ecoding) for diffusion language models, which thresholds a full-span logit-margin aggregate against a smooth, progress-dependent schedule (linear, cosine, exponential), enabling stable, architecture-agnostic stopping without retraining.
2. **Strong efficiency at near-parity quality.** On instruction-tuned models, *SchED* attains average speedups of approximately $4\times$ while maintaining baseline performance. On base models, it delivers consistent acceleration across a spectrum of quality–efficiency trade-offs. In addition to near-parity averages, we observe that *SchED* can yield more correct outputs than vanilla diffusion decoding by terminating at a point of maximum prediction stability, thereby avoiding late-step degradation (Appendix D).
3. **Principled quality–speed trade-off metric.** We introduce *Quality–Penalized Speed (QPS)* (Eq. 16), a conservative metric that jointly accounts for acceleration and quality degradation by explicitly penalizing accuracy drops. Under this metric, *SchED* consistently outperforms prior confidence-based early-exit methods on both base and instruction-tuned dLLMs (Table 3).
4. **Mechanistic explanation via entropy-based analysis.** Analyzing entropy trajectories over diffusion steps show that instruction tuning accelerates confidence stabilization, aligning with progress-aware thresholds and enabling early exit without quality loss (Fig. 2).

SchED builds on the vanilla dLLM denoising process by enforcing a smooth, progress-dependent confidence threshold over the generated span, exploiting the bidirectional, parallel strengths of dLLMs to achieve end-to-end efficiency without further retraining (Zou et al., 2023; Nie et al., 2025; Ye et al., 2025; Gong et al., 2024a). Code is publicly available¹.

2 Related Work

Diffusion LLMs. Recent work establishes diffusion language models (dLLMs) as a viable alternative to autoregressive (AR) models, training either from scratch or by adapting AR checkpoints. LLaDA trains a bidirectional Transformer with a forward masking process and a reverse denoising chain, demonstrating strong pretraining and SFT performance (Nie et al., 2025). Several lines adapt AR models into diffusion ones for fair, scalable comparisons across sizes and tasks (Gong et al., 2024b), while Dream-7B advances recipe and refinement design and releases base/instruct variants (Ye et al., 2025). Architectural advancements have also been integrated, with models like LLaDA-MoE demonstrating the successful application of sparse Mixture-of-Experts to dLLMs, achieving competitive performance with significantly reduced active parameters (Huang et al., 2025). The capabilities of dLLMs are also expanding to long-context scenarios; work like UltraLLaDA shows that techniques such as Rotary Positional Embedding (RoPE) scaling can be adapted to extend context windows to 128K tokens (Wang et al., 2025c). Seed Diffusion emphasizes high-throughput inference for code models (Song et al., 2025). Moreover, formal studies on the scaling laws of masked diffusion models have further established their viability, demonstrating a scaling rate comparable to autoregressive models (Li et al., 2024). Survey and perspective papers summarize the landscape, strengths (parallelism, infilling, global context), and open questions (Zou et al., 2023).

Sampling efficiency and early exit. A central bottleneck for dLLMs is the number of refinement steps. Prophet exploits the empirical observation that answers often stabilize well before the final step and proposes a training-free early-exit rule based on top-2 logit gaps (Pengxiang et al., 2025). Complementary dynamic

policies include SlowFast Sampling, which alternates exploratory vs. accelerated phases guided by certainty/convergence/position principles and can combine with caching (Wei et al., 2025); Duo transfers continuous-diffusion techniques (curriculum learning, consistency distillation) to discrete diffusion to enable few-step sampling (Sahoo et al., 2025); and DLM-One studies one-step generation via score distillation in continuous spaces (Chen et al., 2025). Orthogonal to reducing step counts, other efforts focus on reducing per-step latency by developing novel caching mechanisms like dKV-Cache and d²Cache, which adapt key-value caching to the bidirectional nature of dLLMs (Yuan et al., 2025b; Wang et al., 2025a). Furthermore, speculative decoding has been adapted for the diffusion framework, using the model’s own predictions to draft and verify multiple tokens in parallel, as seen in Spiffy and Self Speculative Decoding (Sohoni et al., 2025; Yuan et al., 2025a). Our work also treats decoding as a *when-to-stop* problem, but (unlike Prophet) we introduce a smooth confidence *schedule* that relaxes thresholds over progress, yielding stable, training-free early exit that is *agnostic to architecture and training*.

3 Methods

In this section, we formalize discrete diffusion language models and introduce *SchED*, our schedule-based early decoding mechanism. We briefly review the forward and reverse masking processes and the standard training objective, then present the confidence aggregation scheme, progress-dependent threshold schedules, and early-exit decoding algorithm.

3.1 Preliminaries: Discrete Diffusion Language Models

Setup and notation. Let \mathcal{V} denote the vocabulary and $[mask]$ a special placeholder token. Given a prompt prefix, generation proceeds over T reverse-diffusion steps $t \in \{1, \dots, T\}$ on sequences $x_t \in (\mathcal{V} \cup \{[mask]\})^L$. At each step, given the current noised sequence x_t and prompt x_{prompt} , the model produces logits

$$L_t = f_\theta(x_{prompt}, x_t, t) \in \mathbb{R}^{L \times |\mathcal{V}|}, \quad (1)$$

and per-position categorical distributions

$$p_{t,i}(\cdot | x_{prompt}, x_t) = \text{softmax}(L_{t,i}) \in \Delta^{|\mathcal{V}|}. \quad (2)$$

¹<https://github.com/amr-mohamedd/SchED>

We write $A \subseteq \{1, \dots, L\}$ for the answer region used later for confidence aggregation, and define normalized diffusion progress $p = t/T$.

Forward (masking) process q . We model discrete diffusion as a progressive *masking process* that corrupts a clean sequence $x_0 \in \mathcal{V}^L$ over T steps. Let $\beta_t \in [0, 1)$ be a step-dependent masking rate and define the per-step transition

$$q(x_t | x_{t-1}) = \prod_{i=1}^L \left[(1 - \beta_t) \delta(x_{t,i} = x_{t-1,i}) + \beta_t \delta(x_{t,i} = [\text{mask}]) \right]. \quad (3)$$

Writing $\bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$ for the token survival probability after t steps, the t -step marginal becomes

$$q(x_t | x_0, t) = \prod_{i=1}^L \left[\bar{\alpha}_t \delta(x_{t,i} = x_{0,i}) + (1 - \bar{\alpha}_t) \delta(x_{t,i} = [\text{mask}]) \right]. \quad (4)$$

Reverse (denoising) process p_θ . Diffusion language models define a learned reverse chain that progressively *unmasks* tokens from x_T (fully masked) to x_0 (fully decoded) by factoring over individual positions:

$$p_\theta(x_{t-1} | x_{\text{prompt}}, x_t) = \prod_{i=1}^L p_\theta(x_{t-1,i} | x_{\text{prompt}}, x_t). \quad (5)$$

Each per-position term is then parameterized as a categorical distribution based on the model’s current output:

$$p_\theta(x_{t-1,i} | x_{\text{prompt}}, x_t) = \text{Cat}(x_{t-1,i}; p_{t,i}(\cdot | x_{\text{prompt}}, x_t)). \quad (6)$$

where $p_{t,i}(\cdot | x_{\text{prompt}}, x_t)$ denotes the model’s categorical prediction at position i given the partially masked sequence x_t and prompt context x_{prompt} at timestep t .

Masked diffusion with partial unmasking (transfer schedule). To control generation granularity, only a subset of masked positions are updated at each step. Let $M_t = \{i : x_{t,i} = [\text{mask}]\}$ and let $\pi_t(i) \in \{0, 1\}$ indicate whether position i is selected for update under a *transfer schedule*. Given the model’s predictive distribution

$p_{t,i}(\cdot | x_{\text{prompt}}, x_t)$ at timestep t , the step operator is

$$x_{t-1,i} = \begin{cases} x_{t,i}, & \text{if } i \notin M_t \text{ or } \pi_t(i) = 0, \\ \hat{x}_{t,i}, & \text{if } i \in M_t \text{ and } \pi_t(i) = 1. \end{cases} \quad (7)$$

where $\hat{x}_{t,i}$ is obtained either deterministically as

$$\hat{x}_{t,i} = \arg \max_v p_{t,i}(v | x_{\text{prompt}}, x_t), \quad (8)$$

or stochastically as

$$\hat{x}_{t,i} \sim \text{Cat}(p_{t,i}(\cdot | x_{\text{prompt}}, x_t)). \quad (9)$$

Block-diffusion variants (e.g., LLaDA) choose π_t over contiguous token blocks, whereas single-block variants (e.g., Dream) often select the entire suffix or a fixed proportion.

Training objective. Diffusion language models are typically trained to invert the forward masking process via a masked-token objective over uniformly sampled timesteps:

$$\mathcal{L}(\theta) = \mathbb{E}_{x_0 \sim \mathcal{D}} \mathbb{E}_{t \sim \mathcal{U}\{1:T\}} \mathbb{E}_{x_t \sim q(\cdot | x_0, t)} \left[- \sum_{i \in M_t} \log p_{t,i}(x_{0,i} | x_{\text{prompt}}, x_t) \right]. \quad (10)$$

Here x_t is obtained by applying the forward kernel at timestep t , and $p_{t,i}(x_{0,i} | x_{\text{prompt}}, x_t)$ is shorthand for the model’s predictive probability $p_\theta(x_{0,i} | x_{\text{prompt}}, x_t, t)$ at position i given the partially noised context. This formulation encourages the model to predict the original tokens at masked sites given partially noised contexts, aligning p_θ with the forward corruption kernel q .

3.2 SchED: Schedule-based Early Decoding

We propose *SchED*, a confidence- and progress-aware early-exit algorithm for diffusion decoding. At each step, *SchED* thresholds the model’s aggregated token confidence against a smooth, nonincreasing function of progress p , i.e., $(\bar{g}_t \geq \tau(p))$. This design builds on the assumption that per-token confidence typically rises as denoising proceeds, allowing generation to terminate precisely when predictions stabilize rather than at a fixed budget. The complete, model-agnostic procedure is summarized in Algorithm 1.

SchED tracks how model confidence evolves throughout denoising and terminates the process once confidence exceeds a progress-dependent threshold, preventing redundant refinement after predictions have stabilized.

Algorithm 1 *SchED*: Schedule-Based Early Decoding for Diffusion LMs

Require: Model M ; prompt x_{prompt} ; gen length L_{gen} ; steps T ; answer region A ; aggregator Agg ; threshold schedule $\tau(p)$; transfer schedule π_t .

Ensure: Completed sequence x

```
1:  $x \leftarrow [x_{\text{prompt}}; [\text{mask}] \times L_{\text{gen}}]$ 
2: for  $t = 1$  to  $T$  do
3:    $L_t \leftarrow M(x)$  ▷ logits
4:    $g_{t,i} \leftarrow L_{t,i}^{(1)} - L_{t,i}^{(2)}$  for  $i \in A$ 
5:    $\bar{g}_t \leftarrow \text{Agg}(\{g_{t,i} : i \in A\})$ 
6:    $p \leftarrow t/T$ 
7:   if  $\bar{g}_t \geq \tau(p)$  then
8:     Fill remaining masks with  $\text{argmax}$ ; return  $x$ 
9:   end if
10:   $M_t \leftarrow \{i : x_i = [\text{mask}]\}$ 
11:   $S_t \leftarrow \{i \in M_t : \pi_t(i) = 1\}$ 
12:   $x_{S_t} \leftarrow \text{arg max}_v L_{t,S_t,v}$ 
13: end for
14: return  $x$ 
```

Confidence measurement. Following the work of Pengxiang et al. (2025), we quantify model confidence using token-level logit margins, which capture how decisively the model prefers its top prediction over alternatives. A region-level confidence score is obtained by aggregating top-2 margins over the *entire answer region* A (i.e., the full model response span):

$$\bar{g}_t = \text{Agg}(\{g_{t,i} : i \in A\}). \quad (11)$$

We use $\text{Agg} = \text{mean}$ by default, so that $\bar{g}_t = \frac{1}{|A|} \sum_{i \in A} g_{t,i}$ in our experiments. Importantly, $g_{t,i}$ uses the *current* logits at step t for every $i \in A$. All thresholds $(\tau_{\text{high}}, \tau_{\text{low}})$ are expressed in *logit units*.

Early-exit trigger. An early exit is triggered when the aggregated confidence surpasses a smooth, progress-dependent threshold schedule $\tau(p)$:

$$\bar{g}_t \geq \tau(p), \quad (12)$$

where $\tau : [0, 1] \rightarrow \mathbb{R}_{\geq 0}$ is a nonincreasing function of the generation progress $p = t/T$. This formulation ensures that the confidence requirement for stopping is highest at the beginning of generation and gradually relaxes as denoising proceeds, allowing the model to terminate decoding once its predictions have stabilized.

Smooth threshold schedules. The threshold function $\tau(p)$ controls when the sampler terminates. We parameterize τ using a pair $(\tau_{\text{high}}, \tau_{\text{low}})$, which specify the initial and final confidence thresholds, and an optional slope parameter $k > 0$. We instantiate the stopping criterion using three distinct

schedule families, each offering a distinct relaxation profile over the diffusion progress. First, the linear schedule provides a constant rate of threshold decay, serving as a straightforward interpolation between the initial and final confidence bounds:

$$\tau_{\text{lin}}(p) = \tau_{\text{high}} + (\tau_{\text{low}} - \tau_{\text{high}})p. \quad (13)$$

Second, we employ a cosine schedule, which facilitates a smooth, sigmoid-like transition by maintaining higher thresholds during the critical early stages of denoising before accelerating relaxation as the process nears completion:

$$\tau_{\text{cos}}(p) = \tau_{\text{low}} + \frac{1}{2}(\tau_{\text{high}} - \tau_{\text{low}}) (1 + \cos(\pi p)). \quad (14)$$

Finally, we introduce an exponential schedule, where the decay trajectory is governed by a curvature parameter $k > 0$.

$$\tau_{\text{exp}}(p) = \tau_{\text{low}} + (\tau_{\text{high}} - \tau_{\text{low}})e^{-kp}. \quad (15)$$

Each of these families defines a continuous, non-increasing trajectory from τ_{high} at the onset ($p = 0$) to τ_{low} at completion ($p = 1$). By parameterizing the exit logic in this manner, *SchED* enables a nuanced trade-off between early-exit aggressiveness and predictive stability, circumventing the brittleness associated with fixed budgets or discrete heuristic rules.

4 Experimental Settings

We evaluate *SchED* across a diverse suite of multiple-choice and long-form generation tasks, assessing its ability to accelerate diffusion decoding while preserving output quality. Specifically, we compare its denoising efficiency against two baselines: (i) standard diffusion sampling without early exit and (ii) Prophet (Pengxiang et al., 2025). Unless otherwise stated, we fix the upper threshold at $(\tau_{\text{high}} = 7.5)$ and, for each schedule family, evaluate two lower-threshold settings: a relaxed $(\tau_{\text{low}} = 0)$ and a more conservative $(\tau_{\text{low}} = 2.5)$, yielding a smooth, monotonic relaxation over diffusion progress in both regimes.

Models. *SchED* is model-agnostic and can be applied to any dLLM without architectural or training modifications. To demonstrate its generality, we evaluate on two representative dLLM families that employ distinct decoding paradigms: *Dream* (Ye et al., 2025), a 7B-parameter model that

performs full-sequence refinement using a single-block decoder, and *LLaDA* (Nie et al., 2025), an 8B-parameter model that adopts a block-diffusion strategy to denoise contiguous token segments. Each model is evaluated in both base and instruction-tuned variants, and we apply the low-confidence remasking strategy across all models.

Benchmarks. We evaluate *SchED* on GPQA (Rein et al., 2023), GSM8K (Cobbe et al., 2021), HellaSwag (Zellers et al., 2019), MMLU (Hendrycks et al., 2020), PIQA (Bisk et al., 2020), and Winogrande (Sakaguchi et al., 2020). For long-form evaluation, we use tasks from the LongBench (Bai et al., 2023), specifically LongBench–HotpotQA (Yang et al., 2018) and LongBench–MultiNews (Fabbri et al., 2019). For translation, we use WMT14 En–Fr (Bojar et al., 2014) and WMT16 En–De (Bojar et al., 2016). Each benchmark includes runs for the baseline, Prophet, and the full set of linear, cosine, and exponential schedules across both model variants.

Evaluation metrics and framework. For multiple-choice (MCQ) benchmarks, *GPQA*, *HellaSwag*, *MMLU*, *PIQA*, and *Winogrande*, we report **accuracy**. For *GSM8K*, we also report **accuracy** based on the generated final answer. For HotpotQA, we report token-level **F1-score**; for MultiNews, we report **ROUGE** (Lin, 2004); and for translation (*WMT14 En–Fr*, *WMT16 En–De*), we report **CHRf**. All evaluations are conducted using the *Language Model Evaluation Harness* (Gao et al., 2024). Complete decoding hyperparameters (step budgets T , generation lengths, shots, and LLaDA block sizes) are provided in Appendix A (Table 4); hardware and total compute are reported in Appendix B.

Efficiency metric for quality–speed trade-offs. Reporting quality and speedup separately can obscure Pareto differences between methods. To jointly capture both aspects, we define the *Quality–Penalized Speed* (QPS) metric:

$$\text{QPS}_\gamma = \text{speedup} \times \left(\frac{\text{score}}{\text{baseline score}} \right)^\gamma, \quad (16)$$

where the exponent $\gamma \geq 1$ controls how strongly quality degradation is penalized. Higher values of γ emphasize fidelity over raw acceleration; in our experiments, we use $\gamma=4$ to provide conservative and interpretable efficiency comparisons.

5 Results

Tables 1 and 2 report results for Dream Base and Dream Instruct, respectively.

Dream Base. With conservative smooth schedules $(\tau_{\text{high}}, \tau_{\text{low}}) = (7.5, 2.5)$, linear, cosine, and $\text{exp-}k=2$, we observe steady $1.04\text{--}1.14\times$ average speedups at near-parity quality, with marginal task gains (e.g., *WMT14 En–Fr*: $\Delta+0.27$) under $\text{exp-}k=2$ (7.5, 2.5). Fast-decaying exponentials (large k) with $\tau_{\text{low}} = 2.5$ remain close to parity, but do not increase the mean speed beyond $\approx 1.1\times$; setting $\tau_{\text{low}} = 0$ yields $2.34\times$ average speed ($\text{exp-}k=16$), with task-specific gains that are most pronounced on *HotpotQA* ($\Delta+7.1$); however, at an overall quality cost of $\approx 2\%$. Prophet offers $\approx 1.07\times$ average speed with near-parity accuracy, providing limited practical benefits relative to smooth schedules.

Dream Instruct. Under the same conservative thresholds, smooth schedules deliver $\approx 3.8\text{--}4.0\times$ average speedups at near parity: table averages cluster around the baseline (e.g., 58.06–58.22 vs. 58.20). Notably, the large- k (fast-decaying) exponential with (7.5, 2.5) attains $4.48\times$, while remaining close to parity (57.59, $\Delta-1\%$). Less conservative settings with $\tau_{\text{low}} = 0$ also preserve translation near baseline with $2.3\text{--}2.8\times$ speed (e.g., cosine and linear). For the large- k exponential, ($k=16$, $\tau_{\text{low}}=0$) delivers the most pronounced speedups while remaining near parity on most benchmarks; however, it exhibits large degradations on *HellaSwag* and *PIQA*.

LLaDA exhibits similar trends to those observed on Dream: conservative schedules provide reliable speedups at near-parity quality on the base model; the instruction-tuned variant yields higher speedups. Fast-decaying settings risk overcommitment that degrades math and long-form. While Prophet provides high speedups, it underperforms on long-form generation. Detailed LLaDA Base and Instruct results are in Tables 7 and 8; additional schedule ablations for Dream and LLaDA appear in Appendix C.

Efficiency results. Table 3 reports *Quality–Penalized Speed* (QPS; $\gamma=4$). On *Dream Base*, smooth schedules concentrate in the $\approx 1.01\text{--}1.12$ range, reflecting the near-parity averages in Table 1 (e.g., 55.02–55.32 vs. baseline 55.31) coupled with modest mean speedups ($\approx 1.04\text{--}1.14\times$). The highest base-model efficiency is achieved by the rapidly decaying exponential,

Method	GPQA	Hella Swag	MMLU	PIQA	Wino grande	GSM8K	Hotpot QA	Multi News	WMT14 En-Fr	WMT16 En-De	Average
Baseline	28.57	79.15	73.39	88.68	75.61	77.10	8.67	21.24	52.99	47.70	55.31
	(×1.00)	(×1.00)	(×1.00)	(×1.00)	(×1.00)	(×1.00)	(×1.00)	(×1.00)	(×1.00)	(×1.00)	(×1.00)
Prophet	28.79	79.15	73.37	88.68	75.61	77.03	8.68	21.21	52.99	47.57	55.31
	(×1.13)	(×1.00)	(×1.03)	(×1.20)	(×1.00)	(×1.07)	(×1.03)	(×1.06)	(×1.11)	(×1.06)	(×1.07)
Cosine (7.5, 0)	28.12	79.15	73.40	88.57	75.61	76.95	8.74	21.28	52.79	45.64	55.02
	(×1.23)	(×1.02)	(×1.09)	(×1.22)	(×1.01)	(×1.20)	(×1.12)	(×1.13)	(×1.23)	(×1.20)	(×1.14)
Cosine (7.5, 2.5)	28.79	79.15	73.37	88.68	75.61	77.03	8.68	21.21	52.99	47.64	55.31
	(×1.11)	(×1.00)	(×1.02)	(×1.14)	(×1.00)	(×1.07)	(×1.02)	(×1.05)	(×1.10)	(×1.05)	(×1.06)
Linear (7.5, 0)	29.02	79.15	73.40	88.57	75.61	76.95	8.66	21.26	52.87	46.43	55.19
	(×1.19)	(×1.00)	(×1.06)	(×1.22)	(×1.00)	(×1.13)	(×1.08)	(×1.10)	(×1.18)	(×1.15)	(×1.11)
Linear (7.5, 2.5)	28.79	79.15	73.37	88.68	75.61	77.03	8.67	21.21	53.00	47.67	55.32
	(×1.10)	(×1.00)	(×1.02)	(×1.03)	(×1.00)	(×1.06)	(×1.02)	(×1.05)	(×1.09)	(×1.04)	(×1.04)
Exp- $k=2$ (7.5, 0)	24.33	79.13	73.50	88.47	75.77	77.63	8.17	20.94	53.16	46.81	54.79
	(×1.19)	(×1.00)	(×1.07)	(×1.27)	(×1.00)	(×1.11)	(×1.07)	(×1.10)	(×1.18)	(×1.13)	(×1.11)
Exp- $k=2$ (7.5, 2.5)	24.33	79.13	73.49	88.63	75.77	77.71	8.16	20.88	53.26	47.60	54.90
	(×1.10)	(×1.00)	(×1.02)	(×1.06)	(×1.00)	(×1.06)	(×1.03)	(×1.05)	(×1.09)	(×1.04)	(×1.04)
Exp- $k=16$ (7.5, 0)	24.78	77.14	73.05	87.92	72.85	62.77	15.75	21.74	52.87	44.75	53.36
	(×2.04)	(×2.50)	(×2.50)	(×2.50)	(×2.50)	(×3.05)	(×2.78)	(×2.75)	(×1.44)	(×1.39)	(×2.34)
Exp- $k=16$ (7.5, 2.5)	23.88	79.13	73.53	88.47	75.77	77.63	8.22	20.92	53.26	47.51	54.83
	(×1.14)	(×1.00)	(×1.18)	(×1.36)	(×1.00)	(×1.08)	(×1.07)	(×1.07)	(×1.13)	(×1.07)	(×1.11)

Table 1: *Dream Base*: benchmark scores with speedups (×). Two thresholds—conservative (7.5, 2.5) and relaxed (7.5, 0). Column best in **bold**, second-best in *italic*.

Exp- $k=16$, (7.5, 0), with **2.03**, driven by a large average speedup (2.34×), and average-score drop of −1.95. Prophet attains 1.07, consistent with its near-baseline averages and limited acceleration.

On *Dream Instruct*, QPS values increase substantially in line with the larger mean speedups observed in Table 2 (≈ 3.8 – $4.0\times$ for conservative smooth schedules) while maintaining average scores close to baseline. The best overall result, **4.30**, is obtained by Exp- $k=16$, (7.5, 2.5), which combines a near-parity average score (57.59) with a pronounced $4.48\times$ mean speedup. Other smooth schedules lie tightly in the 3.78–3.95 band. Prophet lags at 2.91, reflecting weaker average scores and failures on long-form despite notable raw speed. Overall, *Sched* consistently yields higher QPS₄ than Prophet; the top settings pair near-parity average performance with substantial step reductions, which explains their leading efficiency.

6 Analysis

To better understand the differences in speedup values between *base* and *instruct* models, and the variation in speedups across benchmarks, we analyze how predictive uncertainty evolves over the reverse-diffusion chain. We track per-token predictive entropy on the generated region at each step t . Using the per-position distributions $p_{t,i}$ at state x_t ,

the mean per-token entropy is

$$\bar{H}_t = \frac{1}{|S|} \sum_{i \in S} \left(- \sum_{v \in \mathcal{V}} p_{t,i}(v | x_{\text{prompt}}, x_t) \log p_{t,i}(v | x_{\text{prompt}}, x_t) \right). \quad (17)$$

where $p_{t,i}(v)$ denotes the model’s step- t predictive distribution at position i , computed from logits $L_{t,i}$ given x_{prompt} and x_t (Eq. 6).

Results are presented in Figure 2. *Dream Instruct* starts the denoising phase with higher entropy on math-heavy prompts (*GPQA* and *GSM8K*) and decays more slowly in the earliest steps than on *MMLU*, *HotpotQA*, and *WMT14*, yet the trajectories converge to similarly low final entropies. Notably, the *instruct* curves show a rapid early drop followed by nearly uniform per-step decreases, a stepwise profile consistent with *Dream*’s masked diffusion decoding, where at each reverse-diffusion step the model updates high-confidence tokens in parallel and thereby reduces uncertainty in approximately regular increments across the chain. *Dream*’s Confidence-Aware Rescheduling of Tokens (CART) training further sharpens this behavior by adaptively rescheduling token-level noise, thereby encouraging more confident predictions at positions with stronger local context, particularly near the prompt. By contrast, *Dream Base* shows more overlapping trajectories across benchmarks and higher residual entropy overall, suggesting less

Method	GPQA	Hella Swag	MMLU	PIQA	Wino grande	GSM8K	Hotpot QA	Multi News	WMT14 En-Fr	WMT16 En-De	Average
Baseline	30.36	79.91	80.69	73.02	85.80	74.66	27.51	24.39	55.82	49.85	58.20
	($\times 1.00$)	($\times 1.00$)	($\times 1.00$)	($\times 1.00$)	($\times 1.00$)	($\times 1.00$)	($\times 1.00$)	($\times 1.00$)	($\times 1.00$)	($\times 1.00$)	($\times 1.00$)
Prophet	28.79	26.69	78.73	72.33	72.58	66.54	27.87	2.77	20.17	15.28	41.18
	($\times 16.45$)	($\times 7.81$)	($\times 1.14$)	($\times 1.01$)	($\times 1.03$)	($\times 1.02$)	($\times 4.51$)	($\times 31.21$)	($\times 23.93$)	($\times 27.87$)	($\times 11.60$)
Cosine (7.5, 0)	30.58	78.54	80.81	73.06	85.96	74.11	27.46	24.39	55.82	49.86	58.06
	($\times 15.96$)	($\times 2.11$)	($\times 1.56$)	($\times 1.14$)	($\times 1.32$)	($\times 1.61$)	($\times 3.24$)	($\times 6.60$)	($\times 2.39$)	($\times 2.81$)	($\times 3.87$)
Cosine (7.5, 2.5)	30.58	79.53	80.75	73.08	85.91	74.27	27.47	24.38	55.82	49.85	58.16
	($\times 15.93$)	($\times 2.08$)	($\times 1.43$)	($\times 1.11$)	($\times 1.23$)	($\times 1.28$)	($\times 3.14$)	($\times 6.59$)	($\times 2.34$)	($\times 2.77$)	($\times 3.79$)
Linear (7.5, 0)	30.58	79.00	80.81	73.14	85.96	74.11	27.48	24.39	55.82	49.86	58.11
	($\times 16.08$)	($\times 2.11$)	($\times 1.56$)	($\times 1.14$)	($\times 1.32$)	($\times 1.49$)	($\times 3.26$)	($\times 3.02$)	($\times 2.38$)	($\times 2.80$)	($\times 3.89$)
Linear (7.5, 2.5)	30.58	79.76	80.76	73.09	85.85	74.27	27.47	24.38	55.82	49.85	58.16
	($\times 16.01$)	($\times 2.08$)	($\times 1.43$)	($\times 1.10$)	($\times 1.22$)	($\times 1.27$)	($\times 3.17$)	($\times 6.59$)	($\times 2.33$)	($\times 2.76$)	($\times 3.79$)
Exp- $k=2$ (7.5, 0)	30.58	79.45	80.77	73.11	86.16	74.35	26.71	24.56	55.76	49.84	58.14
	($\times 16.32$)	($\times 2.14$)	($\times 1.62$)	($\times 1.21$)	($\times 1.54$)	($\times 1.67$)	($\times 3.47$)	($\times 6.66$)	($\times 2.40$)	($\times 2.76$)	($\times 3.97$)
Exp- $k=2$ (7.5, 2.5)	30.13	80.21	80.69	73.25	86.24	74.51	26.91	24.60	55.76	49.84	58.21
	($\times 16.18$)	($\times 2.10$)	($\times 1.54$)	($\times 1.11$)	($\times 1.23$)	($\times 1.53$)	($\times 3.23$)	($\times 6.57$)	($\times 2.33$)	($\times 2.71$)	($\times 3.85$)
Exp- $k=16$ (7.5, 0)	29.46	32.98	79.31	58.81	84.55	73.80	26.20	23.99	53.53	49.72	51.12
	($\times 18.64$)	($\times 3.73$)	($\times 5.00$)	($\times 2.50$)	($\times 5.00$)	($\times 5.00$)	($\times 5.79$)	($\times 17.38$)	($\times 5.56$)	($\times 3.17$)	($\times 5.44$)
Exp- $k=16$ (7.5, 2.5)	30.58	79.68	79.31	72.78	84.49	72.61	26.58	24.26	55.76	49.84	57.59
	($\times 17.16$)	($\times 2.13$)	($\times 2.50$)	($\times 1.55$)	($\times 2.29$)	($\times 2.10$)	($\times 4.33$)	($\times 7.21$)	($\times 2.38$)	($\times 2.74$)	($\times 4.48$)

Table 2: *Dream Instruct*: benchmark scores with speedups (\times). Two thresholds—conservative (7.5, 2.5) and relaxed (7.5, 0). Column best in **bold**, second-best in *italic*.

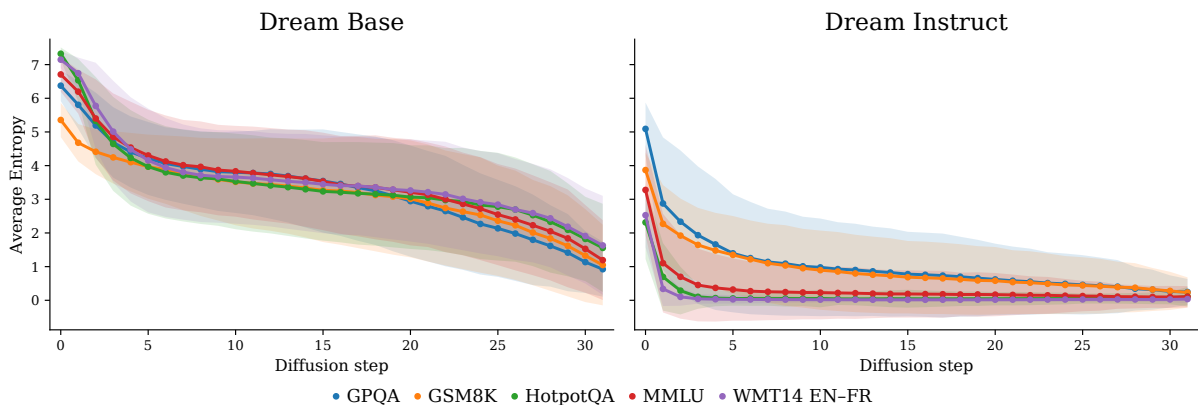


Figure 2: Mean predictive entropy across diffusion steps for five input types (*GPQA*, *GSM8K*, *HotpotQA*, *MMLU*, *WMT14 EN \rightarrow FR*). Curves show per-token entropy (*nats/token*) averaged over evaluation samples; shaded bands denote one standard deviation across samples. Both *Dream Base* and *Dream Instruct* exhibit monotonically decreasing entropy as denoising progresses.

decisive posteriors for QA-style inputs. This aligns with the smaller speedups obtained by schedule-based early exit on Base models (Tables 1–2).

7 Discussion

The robustness of *SchED* across diverse tasks suggests that progress-aware scheduling effectively captures the non-linear stabilization of dLLM posteriors. By framing inference as a dynamic stopping-time problem rather than a fixed-budget process, we align computation with model convergence rather than a conservative step horizon.

Instruction tuning is a primary catalyst for this efficiency. Our entropy analysis (Section 6) reveals that instruction-tuned models collapse uncertainty

over the answer span more rapidly than their base counterparts. This accelerated decay into sharp posterior distributions allows *SchED* to trigger early exits earlier in the denoising chain without sacrificing fidelity, explaining the $4\times$ speedup gap observed in instruction-tuned variants.

Two design elements are critical to this stability. First, full-span aggregation over the response region prevents premature termination from localized confidence spikes—a failure mode common in prior methods like Prophet. Second, smooth, non-increasing schedules dampen transient logit fluctuations, providing a continuous control knob for the quality–speed trade-off.

Notably, *SchED* does not merely trade quality for

Method	Dream Base	Dream Instruct
Cosine (7.5, 0)	1.12	3.83
Cosine (7.5, 2.5)	1.06	3.78
Linear (7.5, 0)	1.10	3.87
Linear (7.5, 2.5)	1.04	3.78
Exp- $k=2$ (7.5, 0)	1.07	3.95
Exp- $k=2$ (7.5, 2.5)	1.01	3.85
Exp- $k=16$ (7.5, 0)	2.03 [†]	3.24
Exp- $k=16$ (7.5, 2.5)	1.07	4.30 [†]
Prophet	1.07	2.91

Table 3: **QPS** ($\gamma=4$) for all *SchED* variants and Prophet. [†] Higher is better.

speed; it can even outperform vanilla decoding on reasoning-intensive inputs. Qualitative cases (Appendix D) show that additional refinement steps in full-horizon diffusion sometimes introduce “over-refinement” errors, such as numerical slips or malformed code. By halting once the aggregate confidence stabilizes, *SchED* preserves high-quality intermediate states that full-horizon decoding may otherwise corrupt with iterative noise. In these instances, progress-aware thresholds provide a Pareto improvement, simultaneously maximizing throughput and correctness.

8 Conclusion

In this work, we introduced *SchED*, a training-free, progress-aware early-exit mechanism that anchors diffusion decoding to sequence-level confidence stabilization. By formulating inference as a *when-to-stop* problem, *SchED* achieves up to $4\times$ speedups, while maintaining, and occasionally surpassing, baseline correctness by mitigating late-stage iterative noise. The method is model-agnostic, integrating seamlessly with existing dLLM architectures without requiring auxiliary training or heuristic step budgets. *SchED* provides a robust and principled framework for efficient dLLM inference, allowing practitioners to precisely match computational effort to model convergence and reliably deploy diffusion models under simultaneous accuracy and throughput constraints.

Limitations

While *SchED* provides a training-free mechanism for accelerating diffusion decoding, our study has some limitations.

SchED shows an explicit quality–efficiency trade-off through the choice of schedule family and hyperparameters (τ_{high} , τ_{low} , k). Conservative schedules (e.g., linear or cosine with higher final

thresholds) yield near-parity quality but moderate speedups, whereas more extreme, rapidly decaying schedules can deliver larger accelerations at the cost of marginal average degradation. In our experiments, this trade-off is resolved by selecting a small set of global configurations per model and regime; in practice, the “right” schedule is task-dependent.

Additionally, *SchED* operates purely at inference time and treats the schedule as an external control signal. We do not investigate joint training of models and schedules, learned aggregation functions beyond averaging over the answer span, or tighter integration with complementary acceleration techniques, such as speculative or cache-based denoising. These directions could further improve robustness and efficiency but are left for future work.

References

- Yushi Bai, Xin Lv, Jiajie Jin, Yidong Zhang, Hongbo Zhang, Beichen Zheng, Yikai Chen, Jian-Guang Qian, Zeyu Chen, Wen-Han Liu, and 1 others. 2023. LongBench: A bilingual, multitask benchmark for long context understanding. *arXiv preprint arXiv:2308.14508*.
- Yonatan Bisk, Rowan Zellers, Ronan Lebras, Chandra Bhagavatula, and Yejin Choi. 2020. PIQA: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 7432–7439.
- Ondřej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Johannes Leveling, Christof Monz, Pavel Pecina, Matt Post, Hervé Saint-Amand, and 1 others. 2014. Findings of the 2014 workshop on statistical machine translation. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 12–58.
- Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, and 1 others. 2016. Findings of the 2016 conference on machine translation (wmt16). In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 131–198.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Tianqi Chen, Shujian Zhang, and Mingyuan Zhou. 2025. DLM-one: Diffusion language models for one-step sequence generation. *arXiv [cs.CL]*.

- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Alexander R Fabbri, Greg Durrett, Daniel Farcas, and Viv Ng. 2019. Multi-News: A large-scale multi-document summarization dataset and abstractive hierarchical model. *arXiv preprint arXiv:1906.01749*.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, and 5 others. 2024. [The language model evaluation harness](#).
- Chen Gong, Hantian Li, Denny Zhou, and Quoc V Le. 2024a. Beyond autoregression: Discrete diffusion for complex reasoning and planning. *arXiv preprint arXiv:2410.14157*.
- Shansan Gong, Shivam Agarwal, Yizhe Zhang, Jiacheng Ye, Lin Zheng, Mukai Li, Chenxin An, Peilin Zhao, Wei Bi, Jiawei Han, Hao Peng, and Lingpeng Kong. 2024b. Scaling diffusion language models via adaptation from autoregressive models. *arXiv [cs.CL]*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Ippolito, Jared Kaplan, Jana Scheurer, Sandeep Datta, Cameron Foote, and 1 others. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
- Jun-Jie Huang, Zixuan Li, Han-Ning Wu, Zixun Lan, Weizhen Qi, Hangbo Bao, Zewen Chi, Wen-Hong Li, Jun-Cheng Chen, Li-Yuan Wang, Wei-Nan Zhang, and Ting Liu. 2025. Llada-moe: A sparse moe diffusion language model. *arXiv preprint arXiv:2509.24389*.
- Xiang Lisa Li, Jack W Rae, Yiding Jiang, Yixiao Li, Hui-Ling Zhen, Grzegorz Górný, Piotr Nawrot, Tim GJ Rudner, Jean-Baptiste Lespiau, Jonathan Richard Schwarz, Arthur Mensch, Roman Ring, Lorenzo Noci, Trevor Cai, Sebastian Borgeaud, Jeff Stanway, Charlie Chen, David Ding, Daniel Autobiography, and 7 others. 2024. Scaling up masked diffusion models on text. *arXiv preprint arXiv:2410.18514*.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Clara Meister, Ayan Pal, Emma Strubell, Mike Rabbat, and Samar Singh. 2024. Fs-dfm: Fast and accurate long text generation with few-step diffusion language models. *arXiv preprint arXiv:2509.20624*.
- Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-Rong Wen, and Chongxuan Li. 2025. Large language diffusion models. *arXiv [cs.CL]*.
- Friedrich Or-EI, Zheng-Hao Liu, Chen-Lin Zhang, Hong-Yi Yuan, and Wei-Qiang Zhang. 2025. Accelerating diffusion llm inference via local determinism propagation. *arXiv preprint arXiv:2510.07081*.
- Ayan Pal, Meet Shah, Emma Strubell, Berkin Gunel, Clara Meister, Ameet Pokle, Yilun Zhou, Mike Rabbat, and Samar Singh. 2025. Mercury: Ultra-fast language models based on diffusion. *arXiv preprint arXiv:2506.17298*.
- Li Pengxiang, Zhou Yefan, Muhtar Dilxat, Yin Lu, Yan Shilin, Shen Li, Liang Yi, Vosoughi Soroush, and Liu Shiwei. 2025. Diffusion language models know the answer before decoding. *arXiv [cs.CL]*.
- Mary Phuong, Chen Gong, Hantian Li, Denny Zhou, and Quoc V Le. 2025. d1: Scaling reasoning in diffusion large language models via reinforcement learning. *arXiv preprint arXiv:2504.12216*.
- David Rein, Stas Raichuk, Orion Selesnick, Armel Fotso, Jarno Textor, Daniel Lindner, Yafang Du, Tsvi August, Dan Gutman, Mor Geva, and 1 others. 2023. GPQA: A graduate-level google-proof Q&A benchmark. *arXiv preprint arXiv:2311.12022*.
- Subham Sekhar Sahoo, Justin Deschenaux, Aaron Gokaslan, Guanghan Wang, Justin Chiu, and Volodymyr Kuleshov. 2025. The diffusion duality. *arXiv [cs.LG]*.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. WinoGrande: An adversarial Winograd schema challenge at scale. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8732–8740.
- Nimit S Sohoni, Chien-Yu Lin, Aniruddha Nrusimha, and Christopher Ré. 2025. Spiffy: Multiplying diffusion llm acceleration via lossless speculative decoding. *arXiv preprint arXiv:2509.18085*.
- Yuxuan Song, Zheng Zhang, Cheng Luo, Pengyang Gao, Fan Xia, Hao Luo, Zheng Li, Yuehang Yang, Hongli Yu, Xingwei Qu, Yuwei Fu, Jing Su, Ge Zhang, Wenhao Huang, Mingxuan Wang, Lin Yan, Xiaoying Jia, Jingjing Liu, Wei-Ying Ma, and 3 others. 2025. Seed diffusion: A large-scale diffusion language model with high-speed inference. *arXiv [cs.CL]*.
- Yixuan Su, Cheng Zhu, Wangchunshu Zhou, Jing Chen, and Graham Neubig. 2025. Don’t settle too early: Self-reflective remasking for diffusion language models. *arXiv preprint arXiv:2509.23653*.
- Chuhan Wang, Hong-Ning Dai, Yi-Chen Zhang, Chun-Wei Tsung, Xin-Hao Li, Wei-Guo Zheng, and Ming-Li Song. 2025a. d²cache: Accelerating diffusion-based llms via dual adaptive caching. *arXiv preprint arXiv:2509.23094*.

- Chuhan Wang, Yi-Chen Zhang, Xin-Hao Li, Hong-Ning Dai, and Ming-Li Song. 2025b. Creditdecoding: Accelerating parallel decoding in diffusion large language models with trace credits. *arXiv preprint arXiv:2510.06133*.
- Ming-Hao Wang, Yun-Zhu Song, Jun-Jie Huang, Zixuan Li, Wen-Hong Li, Jia-Jun Zhang, and Cheng-Chung Hsu. 2025c. Ultrallada: Scaling the context length to 128k for diffusion large language models. *arXiv preprint arXiv:2510.10481*.
- Qingyan Wei, Yaojie Zhang, Zhiyuan Liu, Dongrui Liu, and Linfeng Zhang. 2025. Accelerating diffusion large language models with SlowFast sampling: The three golden principles. *arXiv [cs.CL]*.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380.
- Jiacheng Ye, Zhihui Xie, Lin Zheng, Jiahui Gao, Zirui Wu, Xin Jiang, Zhenguo Li, and Lingpeng Kong. 2025. Dream 7B: Diffusion large language models. *arXiv [cs.CL]*.
- Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. 2024. A survey on multimodal large language models. *National Science Review*, 11(12):nwae403.
- Hong-Yi Yuan, Zheng-Hao Liu, Chen-Lin Zhang, and Wei-Qiang Zhang. 2025a. Self speculative decoding for diffusion large language models. *arXiv preprint arXiv:2510.04147*.
- Hongyi Yuan, Zheng-Hao Liu, Chen-Lin Zhang, Xuran Wang, Yu-Hong-Vey Wong, and Jia-Wei Liu. 2025b. dkv-cache: The cache for diffusion language models. *arXiv preprint arXiv:2505.15781*.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. HellaSwag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800.
- Chen-Lin Zhang, Zheng-Hao Liu, Hong-Yi Yuan, and Wei-Qiang Zhang. 2025a. Finish first, perfect later: Test-time token-level cross-validation for diffusion large language models. *arXiv preprint arXiv:2510.05090*.
- Lingzhe Zhang, Liancheng Fang, Chiming Duan, Minghua He, Leyi Pan, Pei Xiao, Shiyu Huang, Yunpeng Zhai, Xuming Hu, Philip S Yu, and 1 others. 2025b. A survey on parallel text generation: From parallel decoding to diffusion language models. *arXiv preprint arXiv:2508.08712*.
- Hao Zou, Zae Myung Kim, and Dongyeop Kang. 2023. A survey of diffusion models in natural language processing. *arXiv [cs.CL]*.

A Generation Configurations

Table 4 summarizes the decoding hyperparameters used across benchmarks, as well as the total number of instances evaluated for each dataset. For each task we specify the maximum number of reverse-diffusion steps T (the step budget over which progress $p=t/T$ is normalized), the *generation length* (the maximum number of tokens in the generated answer region A), the number of in-context examples (*shots*), and—for the block-diffusion variant *LLaDA*—the block size used by the transfer schedule π_t (Dream uses single-block/suffix refinement and thus does not require a block-size setting).

Short multiple-choice (MCQ) tasks are evaluated with compact budgets ($T=5$) and very short generation lengths, whereas math, translation, and long-form tasks employ substantially larger step budgets (256–512) and longer generation lengths, making them more sensitive to fast-decaying early-exit schedules. All MCQ benchmarks (MMLU, HellaSwag, PIQA, Winogrande) are evaluated in *generative* mode with full decoding of the answer region, *not* via likelihood/ranking of options; the model must produce the final answer tokens in A , and accuracy is computed from the generated outputs. For *LLaDA*, we use a block size of 5 tokens on short MCQ benchmarks and 32 tokens on tasks that require generation lengths of 32 tokens or more.

Benchmark	N. of Samples	Max Steps T	Generation Length	Shots	LLaDA Block-size
MMLU (generative)	14,042	5	5	5	5
HellaSwag (generative)	10,042	5	5	5	5
PIQA (generative)	1,838	5	5	5	5
Winogrande (generative)	1,267	5	5	5	5
GPQA	448	128	128	8	32
GSM8K	1,319	256	256	8	32
WMT14 En→Fr	3,003	256	256	5	32
WMT16 En→De	2,998	256	256	5	32
LongBench MultiNews	200	512	512	0	32
LongBench HotpotQA	200	32	32	0	32

Table 4: Decoding hyperparameters and instance counts per benchmark. *Generation length* denotes the maximum number of tokens in the generated answer region A . *LLaDA Block* applies only to the block-diffusion decoder; Dream uses a single-block transfer schedule.

B Hardware and Compute Budget

All experiments were conducted on $8 \times$ AMD MI210, with a total compute allocation of 52 GPU-days.

Method	GPQA	HellaSwag	MMLU	PIQA	Winogrande	GSM8K	HotpotQA	MultiNews	WMT14 En-Fr	WMT16 En-De	Average
Baseline	28.57 ($\times 1.00$)	79.15 ($\times 1.00$)	73.39 ($\times 1.00$)	88.68 ($\times 1.00$)	75.61 ($\times 1.00$)	77.10 ($\times 1.00$)	8.67 ($\times 1.00$)	21.24 ($\times 1.00$)	52.99 ($\times 1.00$)	47.70 ($\times 1.00$)	55.31 ($\times 1.00$)
Prophet	28.79 ($\times 1.13$)	79.15 ($\times 1.00$)	73.37 ($\times 1.03$)	88.68 ($\times 1.20$)	75.61 ($\times 1.00$)	77.03 ($\times 1.07$)	8.68 ($\times 1.03$)	21.21 ($\times 1.06$)	52.99 ($\times 1.11$)	47.57 ($\times 1.06$)	55.31 ($\times 1.07$)
Cosine (7.5, 2.5)	28.79 ($\times 1.11$)	79.15 ($\times 1.00$)	73.37 ($\times 1.02$)	88.68 ($\times 1.14$)	75.61 ($\times 1.00$)	77.03 ($\times 1.07$)	8.68 ($\times 1.02$)	21.21 ($\times 1.05$)	52.99 ($\times 1.10$)	47.64 ($\times 1.05$)	55.31 ($\times 1.06$)
Cosine (7.5, 0)	28.12 ($\times 1.23$)	79.15 ($\times 1.02$)	73.40 ($\times 1.09$)	88.57 ($\times 1.22$)	75.61 ($\times 1.01$)	76.95 ($\times 1.20$)	8.74 ($\times 1.12$)	21.28 ($\times 1.13$)	52.79 ($\times 1.23$)	45.64 ($\times 1.20$)	55.02 ($\times 1.14$)
Linear (7.5, 2.5)	28.79 ($\times 1.10$)	79.15 ($\times 1.00$)	73.37 ($\times 1.02$)	88.68 ($\times 1.03$)	75.61 ($\times 1.00$)	77.03 ($\times 1.06$)	8.67 ($\times 1.02$)	21.21 ($\times 1.05$)	53.00 ($\times 1.09$)	47.67 ($\times 1.04$)	55.32 ($\times 1.04$)
Linear (7.5, 0)	29.02 ($\times 1.19$)	79.15 ($\times 1.00$)	73.40 ($\times 1.06$)	88.57 ($\times 1.22$)	75.61 ($\times 1.00$)	76.95 ($\times 1.13$)	8.66 ($\times 1.08$)	21.26 ($\times 1.10$)	52.87 ($\times 1.18$)	46.43 ($\times 1.15$)	55.19 ($\times 1.11$)
Exp- $k=2$ (7.5, 2.5)	24.33 ($\times 1.10$)	79.13 ($\times 1.00$)	73.49 ($\times 1.02$)	88.63 ($\times 1.06$)	75.77 ($\times 1.00$)	77.71 ($\times 1.06$)	8.16 ($\times 1.03$)	20.88 ($\times 1.05$)	53.26 ($\times 1.09$)	47.60 ($\times 1.04$)	54.90 ($\times 1.04$)
Exp- $k=2$ (7.5, 0)	24.33 ($\times 1.19$)	79.13 ($\times 1.00$)	73.50 ($\times 1.07$)	88.47 ($\times 1.27$)	75.77 ($\times 1.00$)	77.63 ($\times 1.11$)	8.17 ($\times 1.07$)	20.94 ($\times 1.10$)	53.16 ($\times 1.18$)	46.81 ($\times 1.13$)	54.79 ($\times 1.11$)
Exp- $k=4$ (7.5, 2.5)	23.66 ($\times 1.13$)	79.13 ($\times 1.00$)	73.49 ($\times 1.06$)	88.63 ($\times 1.20$)	75.77 ($\times 1.00$)	77.63 ($\times 1.07$)	8.16 ($\times 1.05$)	20.90 ($\times 1.06$)	53.24 ($\times 1.12$)	47.51 ($\times 1.05$)	54.81 ($\times 1.07$)
Exp- $k=4$ (7.5, 0)	23.44 ($\times 1.33$)	79.13 ($\times 1.15$)	73.53 ($\times 1.51$)	88.47 ($\times 1.62$)	75.77 ($\times 1.31$)	77.33 ($\times 1.47$)	8.89 ($\times 1.35$)	21.03 ($\times 1.30$)	52.97 ($\times 1.32$)	44.95 ($\times 1.30$)	54.55 ($\times 1.37$)
Exp- $k=8$ (7.5, 2.5)	23.66 ($\times 1.14$)	79.13 ($\times 1.00$)	73.53 ($\times 1.15$)	88.52 ($\times 1.29$)	75.77 ($\times 1.00$)	77.63 ($\times 1.08$)	8.22 ($\times 1.07$)	20.92 ($\times 1.07$)	53.26 ($\times 1.13$)	47.51 ($\times 1.06$)	54.81 ($\times 1.10$)
Exp- $k=8$ (7.5, 0)	24.78 ($\times 1.60$)	78.99 ($\times 1.56$)	73.33 ($\times 2.39$)	87.98 ($\times 2.42$)	73.56 ($\times 2.39$)	72.71 ($\times 2.20$)	14.02 ($\times 2.17$)	21.57 ($\times 1.91$)	52.91 ($\times 1.40$)	44.78 ($\times 1.37$)	54.46 ($\times 1.94$)
Exp- $k=16$ (7.5, 2.5)	23.88 ($\times 1.14$)	79.13 ($\times 1.00$)	73.53 ($\times 1.18$)	88.47 ($\times 1.36$)	75.77 ($\times 1.00$)	79.83 ($\times 2.12$)	8.22 ($\times 1.07$)	20.92 ($\times 1.07$)	53.26 ($\times 1.13$)	47.51 ($\times 1.07$)	54.83 ($\times 1.11$)
Exp- $k=16$ (7.5, 0)	24.78 ($\times 2.04$)	77.14 ($\times 2.50$)	73.05 ($\times 2.50$)	87.92 ($\times 2.50$)	72.85 ($\times 2.50$)	62.77 ($\times 3.05$)	15.75 ($\times 2.78$)	21.74 ($\times 2.75$)	52.87 ($\times 1.44$)	44.75 ($\times 1.39$)	53.36 ($\times 2.34$)

Table 5: **Dream Base**. Highest accuracy per column in **bold**; second-highest in *italics*. The rightmost column reports mean accuracy with mean speedup in parentheses.

Method	GPQA	HellaSwag	MMLU	PIQA	Winogrande	GSM8K	HotpotQA	MultiNews	WMT14 En-Fr	WMT16 En-De	Average
Baseline	30.36 ($\times 1.00$)	79.91 ($\times 1.00$)	80.69 ($\times 1.00$)	73.02 ($\times 1.00$)	85.80 ($\times 1.00$)	74.66 ($\times 1.00$)	27.51 ($\times 1.00$)	24.39 ($\times 1.00$)	55.82 ($\times 1.00$)	49.85 ($\times 1.00$)	58.20 ($\times 1.00$)
Prophet	28.79 ($\times 16.45$)	26.69 ($\times 7.81$)	78.73 ($\times 1.14$)	72.33 ($\times 1.01$)	72.58 ($\times 1.03$)	66.54 ($\times 1.02$)	27.87 ($\times 4.51$)	2.77 ($\times 31.21$)	20.17 ($\times 23.93$)	15.28 ($\times 27.87$)	41.18 ($\times 11.60$)
Cosine (7.5, 2.5)	30.58 ($\times 15.93$)	79.53 ($\times 2.08$)	80.75 ($\times 1.43$)	73.08 ($\times 1.11$)	85.91 ($\times 1.23$)	74.27 ($\times 1.28$)	27.47 ($\times 3.14$)	24.38 ($\times 6.59$)	55.82 ($\times 2.34$)	49.85 ($\times 2.77$)	58.16 ($\times 3.79$)
Cosine (7.5, 0)	30.58 ($\times 15.96$)	78.54 ($\times 2.11$)	80.81 ($\times 1.56$)	73.06 ($\times 1.14$)	85.96 ($\times 1.32$)	74.11 ($\times 1.61$)	27.46 ($\times 3.24$)	24.39 ($\times 6.60$)	55.82 ($\times 2.39$)	49.86 ($\times 2.81$)	58.06 ($\times 3.87$)
Linear (7.5, 2.5)	30.58 ($\times 16.01$)	79.76 ($\times 2.08$)	80.76 ($\times 1.43$)	73.09 ($\times 1.10$)	85.85 ($\times 1.22$)	74.27 ($\times 1.27$)	27.47 ($\times 3.17$)	24.38 ($\times 6.59$)	55.82 ($\times 2.33$)	49.85 ($\times 2.76$)	58.16 ($\times 3.79$)
Linear (7.5, 0)	30.58 ($\times 16.08$)	79.00 ($\times 2.11$)	80.81 ($\times 1.56$)	73.14 ($\times 1.14$)	85.96 ($\times 1.32$)	74.11 ($\times 1.49$)	27.48 ($\times 3.26$)	24.39 ($\times 3.02$)	55.82 ($\times 2.38$)	49.86 ($\times 2.80$)	58.11 ($\times 3.89$)
Exp- $k=2$ (7.5, 2.5)	30.13 ($\times 16.18$)	80.21 ($\times 2.10$)	80.69 ($\times 1.54$)	73.25 ($\times 1.11$)	86.24 ($\times 1.23$)	74.51 ($\times 1.53$)	26.91 ($\times 3.23$)	24.60 ($\times 6.57$)	55.76 ($\times 2.33$)	49.84 ($\times 2.71$)	58.21 ($\times 3.85$)
Exp- $k=2$ (7.5, 0)	30.58 ($\times 16.32$)	79.45 ($\times 2.14$)	80.77 ($\times 1.62$)	73.11 ($\times 1.21$)	86.16 ($\times 1.54$)	74.35 ($\times 1.67$)	26.71 ($\times 3.47$)	24.56 ($\times 6.66$)	55.76 ($\times 2.40$)	49.84 ($\times 2.76$)	58.14 ($\times 3.97$)
Exp- $k=4$ (7.5, 2.5)	30.58 ($\times 16.38$)	80.14 ($\times 2.12$)	80.78 ($\times 1.70$)	73.10 ($\times 1.14$)	86.29 ($\times 2.01$)	74.27 ($\times 1.49$)	26.86 ($\times 3.47$)	24.57 ($\times 6.67$)	55.76 ($\times 2.36$)	49.84 ($\times 2.73$)	58.22 ($\times 3.97$)
Exp- $k=4$ (7.5, 0)	30.58 ($\times 16.62$)	71.27 ($\times 2.25$)	79.31 ($\times 2.50$)	70.57 ($\times 1.64$)	84.82 ($\times 2.51$)	74.66 ($\times 2.54$)	26.58 ($\times 4.03$)	24.03 ($\times 5.59$)	55.78 ($\times 2.94$)	49.84 ($\times 2.85$)	56.57 ($\times 4.44$)
Exp- $k=8$ (7.5, 2.5)	30.80 ($\times 16.69$)	79.83 ($\times 2.12$)	79.48 ($\times 2.48$)	73.53 ($\times 1.24$)	86.13 ($\times 2.37$)	73.64 ($\times 1.72$)	26.65 ($\times 3.82$)	24.29 ($\times 5.50$)	55.76 ($\times 2.37$)	49.84 ($\times 2.74$)	57.97 ($\times 4.23$)
Exp- $k=8$ (7.5, 0)	30.36 ($\times 17.03$)	43.67 ($\times 2.94$)	79.31 ($\times 4.21$)	58.86 ($\times 1.64$)	84.55 ($\times 4.76$)	74.03 ($\times 4.64$)	26.28 ($\times 4.77$)	20.86 ($\times 9.53$)	55.27 ($\times 4.15$)	49.86 ($\times 3.05$)	52.48 ($\times 4.94$)
Exp- $k=16$ (7.5, 2.5)	30.58 ($\times 17.16$)	79.68 ($\times 2.13$)	79.31 ($\times 2.50$)	72.78 ($\times 1.55$)	84.49 ($\times 2.29$)	72.61 ($\times 2.10$)	26.58 ($\times 4.33$)	24.26 ($\times 7.21$)	55.76 ($\times 2.38$)	49.84 ($\times 2.74$)	57.59 ($\times 4.48$)
Exp- $k=16$ (7.5, 0)	29.46 ($\times 18.64$)	32.98 ($\times 3.73$)	79.31 ($\times 5.00$)	58.81 ($\times 2.50$)	84.55 ($\times 5.00$)	73.80 ($\times 5.00$)	26.20 ($\times 5.79$)	23.99 ($\times 17.38$)	53.53 ($\times 5.56$)	49.72 ($\times 3.17$)	51.12 ($\times 5.44$)

Table 6: **Dream Instruct**. Highest accuracy per column in **bold**; second-highest in *italics*. The rightmost column reports mean accuracy with mean speedup in parentheses.

C Additional *Sched* Variants and *LLaDA* Results

This appendix primarily details the full *LLaDA* results—both *Base* and *Instruct*—under the same evaluation protocol and schedule families as the main paper (Tables 7–8). In addition, we report one new set of variants for *Dream*: *intermediate-curvature exponential* schedules with $k \in \{4, 8\}$, each evaluated under the conservative threshold $(\tau_{\text{high}}, \tau_{\text{low}}) = (7.5, 2.5)$ and the relaxed $(7.5, 0)$ (Tables 5–6). All other schedules (linear, cosine, and exponential with $k \in \{2, 16\}$) are already covered in the main text.

Dream. On *Dream Instruct*, the intermediate exponentials are highly competitive with the best smooth schedules. In particular, Exp- $k=4$ with $(7.5, 2.5)$ attains the highest overall average score 58.22 at an average speedup of $\times 3.97$ (Table 6). Exp- $k=8$ with $(7.5, 2.5)$ matches the PIQA peak (73.53) while preserving translation near baseline (WMT14 En-Fr 55.76, WMT16 En-De 49.84) and pushing the mean speedup to $\times 4.23$. At the task level, the intermediate schedules frequently sit on or near the column bests: e.g., Exp- $k=8$ $(7.5, 2.5)$ yields the top GPQA (30.80) and shares the PIQA peak, while Exp- $k=4$ $(7.5, 2.5)$ is within rounding of the HellaSwag best (80.14 vs. 80.21 for Exp- $k=2$). As intended, their average speedups typically fall between the relaxed $k=2$ exponential (e.g., 58.14 at $\times 3.97$ for $(7.5, 0)$) and the high-curvature $k=16$ exponential (57.59 at $\times 4.48$ for $(7.5, 2.5)$).

On *Dream Base* (Table 5), the same pattern holds at a smaller scale. Exp- $k=8$ $(7.5, 2.5)$ reaches the best MMLU (73.53) and ties the top WMT14 En-Fr (53.26) while delivering a mean $\times 1.10$ speedup; Exp- $k=4$ $(7.5, 2.5)$ is similarly competitive on WMT14 (53.24) and PIQA (second-best 88.63) with $\times 1.07$ speed. These sit neatly between the relaxed $k=2$ exponential (mean $\times 1.11$) and the high-curvature $k=16$ exponential (mean $\times 1.11$ under $(7.5, 2.5)$, or much larger $\times 2.34$ under $(7.5, 0)$ with the expected quality trade-off).

Method	GPQA	HellaSwag	MMLU	PIQA	Winogrande	GSM8K	HotpotQA	MultiNews	WMT14 En-Fr	WMT16 En-De	Average
Baseline	26.12 ($\times 1.00$)	86.20 ($\times 1.00$)	58.16 ($\times 1.00$)	<i>81.72</i> ($\times 1.00$)	77.98 ($\times 1.00$)	47.76 ($\times 1.00$)	9.08 ($\times 1.00$)	23.85 ($\times 1.00$)	62.39 ($\times 1.00$)	<i>56.67</i> ($\times 1.00$)	52.99 ($\times 1.00$)
Prophet	31.03 ($\times 15.32$)	<i>86.17</i> ($\times 1.24$)	<i>58.00</i> ($\times 1.24$)	<i>81.56</i> ($\times 2.05$)	77.98 ($\times 1.25$)	16.07 ($\times 7.84$)	<i>14.94</i> ($\times 3.19$)	18.83 ($\times 32.48$)	<i>49.58</i> ($\times 17.53$)	<i>45.15</i> ($\times 17.98$)	<i>47.93</i> ($\times 10.01$)
Cosine (7.5, 2.5)	27.23 ($\times 1.97$)	<i>86.16</i> ($\times 1.17$)	<i>58.00</i> ($\times 1.24$)	81.77 ($\times 2.04$)	77.98 ($\times 1.23$)	39.12 ($\times 1.97$)	<i>9.78</i> ($\times 1.60$)	25.00 ($\times 1.80$)	62.41 ($\times 1.87$)	56.68 ($\times 1.87$)	52.41 ($\times 1.68$)
Cosine (7.5, 0)	27.90 ($\times 2.24$)	<i>86.08</i> ($\times 1.49$)	<i>58.00</i> ($\times 1.24$)	<i>81.50</i> ($\times 2.12$)	77.98 ($\times 1.54$)	42.68 ($\times 2.24$)	11.12 ($\times 1.88$)	25.23 ($\times 2.07$)	62.41 ($\times 2.10$)	56.66 ($\times 2.10$)	52.96 ($\times 1.90$)
Exp- $k=16$ (7.5, 2.5)	<i>30.80</i> ($\times 4.47$)	<i>85.98</i> ($\times 1.67$)	<i>58.00</i> ($\times 1.24$)	80.30 ($\times 4.69$)	77.98 ($\times 1.65$)	30.40 ($\times 3.75$)	12.93 ($\times 2.84$)	24.77 ($\times 3.38$)	62.19 ($\times 3.09$)	56.40 ($\times 3.04$)	51.98 ($\times 2.98$)
Exp- $k=16$ (7.5, 0)	<i>28.79</i> ($\times 8.67$)	<i>86.01</i> ($\times 5.00$)	<i>58.00</i> ($\times 1.24$)	79.98 ($\times 5.00$)	77.66 ($\times 5.00$)	4.93 ($\times 8.58$)	20.98 ($\times 7.13$)	20.74 ($\times 8.55$)	51.15 ($\times 7.93$)	47.03 ($\times 7.79$)	47.53 ($\times 6.49$)
Exp- $k=2$ (7.5, 2.5)	26.34 ($\times 2.24$)	<i>86.16</i> ($\times 1.07$)	<i>58.00</i> ($\times 1.24$)	<i>81.56</i> ($\times 2.12$)	77.98 ($\times 1.14$)	42.38 ($\times 2.19$)	9.83 ($\times 1.65$)	25.11 ($\times 1.92$)	62.41 ($\times 1.99$)	56.68 ($\times 2.00$)	52.65 ($\times 1.76$)
Exp- $k=2$ (7.5, 0)	29.02 ($\times 2.85$)	<i>86.03</i> ($\times 1.58$)	<i>58.00</i> ($\times 1.24$)	<i>81.23</i> ($\times 2.45$)	77.98 ($\times 1.64$)	43.52 ($\times 2.75$)	11.65 ($\times 2.20$)	25.30 ($\times 2.47$)	62.38 ($\times 2.46$)	56.62 ($\times 2.45$)	53.17 ($\times 2.21$)
Exp- $k=4$ (7.5, 2.5)	28.79 ($\times 2.86$)	<i>86.06</i> ($\times 1.41$)	<i>58.00</i> ($\times 1.24$)	<i>81.18</i> ($\times 2.42$)	77.98 ($\times 1.43$)	<i>43.75</i> ($\times 2.73$)	11.57 ($\times 2.08$)	25.37 ($\times 2.40$)	62.40 ($\times 2.39$)	56.64 ($\times 2.39$)	53.17 ($\times 2.13$)
Exp- $k=4$ (7.5, 0)	27.46 ($\times 4.09$)	<i>85.96</i> ($\times 2.34$)	<i>58.00</i> ($\times 1.24$)	<i>81.28</i> ($\times 2.48$)	78.53 ($\times 3.17$)	31.01 ($\times 3.72$)	13.59 ($\times 3.17$)	24.72 ($\times 3.52$)	61.82 ($\times 3.37$)	56.08 ($\times 3.33$)	51.85 ($\times 3.05$)
Exp- $k=8$ (7.5, 2.5)	28.57 ($\times 3.70$)	<i>86.05</i> ($\times 1.62$)	<i>58.00</i> ($\times 1.24$)	<i>81.28</i> ($\times 2.95$)	77.98 ($\times 1.58$)	37.23 ($\times 3.32$)	12.50 ($\times 2.57$)	25.10 ($\times 3.00$)	62.30 ($\times 2.81$)	56.53 ($\times 2.80$)	52.55 ($\times 2.55$)
Exp- $k=8$ (7.5, 0)	28.35 ($\times 5.94$)	<i>85.95</i> ($\times 3.37$)	<i>58.00</i> ($\times 1.24$)	79.98 ($\times 5.00$)	78.69 ($\times 3.99$)	17.36 ($\times 5.46$)	<i>17.04</i> ($\times 5.74$)	23.45 ($\times 9.53$)	58.77 ($\times 5.00$)	53.58 ($\times 4.95$)	50.12 ($\times 4.49$)
Linear (7.5, 2.5)	28.12 ($\times 1.98$)	<i>86.16</i> ($\times 1.07$)	<i>58.00</i> ($\times 1.24$)	<i>81.61</i> ($\times 2.07$)	77.98 ($\times 1.14$)	39.12 ($\times 1.96$)	9.68 ($\times 1.53$)	24.98 ($\times 1.76$)	62.41 ($\times 1.85$)	56.68 ($\times 1.85$)	52.47 ($\times 1.65$)
Linear (7.5, 0)	28.57 ($\times 2.35$)	<i>86.08</i> ($\times 1.41$)	<i>58.00</i> ($\times 1.24$)	<i>81.61</i> ($\times 2.14$)	77.98 ($\times 1.42$)	42.84 ($\times 2.32$)	11.05 ($\times 1.86$)	25.25 ($\times 2.09$)	62.41 ($\times 2.13$)	56.66 ($\times 2.13$)	53.05 ($\times 1.91$)

Table 7: **LLaDA Base** with intermediate-curvature exponentials ($k \in \{4, 8\}$) and explicit thresholds. Highest accuracy per column in **bold**; second-highest in *italics*. The rightmost column reports mean accuracy with mean speedup in parentheses.

Method	GPQA	HellaSwag	MMLU	PIQA	Winogrande	GSM8K	HotpotQA	MultiNews	WMT14 En-Fr	WMT16 En-De	Average
Baseline	24.33 ($\times 1.00$)	74.28 ($\times 1.00$)	63.19 ($\times 1.00$)	82.86 ($\times 1.00$)	76.72 ($\times 1.00$)	51.93 ($\times 1.00$)	10.65 ($\times 1.00$)	26.79 ($\times 1.00$)	60.80 ($\times 1.00$)	54.26 ($\times 1.00$)	52.28 ($\times 1.00$)
Prophet	23.44 ($\times 5.97$)	<i>74.74</i> ($\times 1.93$)	63.91 ($\times 1.63$)	82.70 ($\times 1.81$)	<i>76.16</i> ($\times 1.45$)	27.60 ($\times 18.62$)	9.37 ($\times 3.31$)	19.25 ($\times 37.98$)	44.32 ($\times 26.05$)	45.71 ($\times 15.26$)	46.15 ($\times 11.52$)
Cosine (7.5, 2.5)	24.78 ($\times 1.89$)	74.69 ($\times 1.92$)	<i>63.86</i> ($\times 1.64$)	82.70 ($\times 1.78$)	<i>76.16</i> ($\times 1.45$)	37.30 ($\times 2.25$)	<i>9.81</i> ($\times 1.76$)	26.80 ($\times 2.35$)	<i>55.39</i> ($\times 24.65$)	<i>47.93</i> ($\times 29.64$)	49.74 ($\times 6.83$)
Cosine (7.5, 0)	25.22 ($\times 2.18$)	<i>74.69</i> ($\times 1.93$)	<i>63.86</i> ($\times 1.64$)	82.70 ($\times 1.78$)	<i>76.16</i> ($\times 1.57$)	36.85 ($\times 2.50$)	9.16 ($\times 2.07$)	26.68 ($\times 2.63$)	55.29 ($\times 24.66$)	47.87 ($\times 29.65$)	49.74 ($\times 7.06$)
Exp- $k=16$ (7.5, 2.5)	24.33 ($\times 4.19$)	74.63 ($\times 2.22$)	<i>63.86</i> ($\times 1.64$)	<i>81.56</i> ($\times 2.82$)	76.09 ($\times 2.10$)	33.59 ($\times 5.74$)	9.18 ($\times 4.10$)	22.37 ($\times 8.57$)	37.47 ($\times 35.10$)	33.91 ($\times 37.22$)	47.48 ($\times 10.87$)
Exp- $k=16$ (7.5, 0)	<i>25.00</i> ($\times 9.26$)	73.51 ($\times 5.00$)	<i>63.86</i> ($\times 1.64$)	79.00 ($\times 5.00$)	73.80 ($\times 5.00$)	27.98 ($\times 11.45$)	7.61 ($\times 9.17$)	16.78 ($\times 17.38$)	33.27 ($\times 44.67$)	29.45 ($\times 45.10$)	45.22 ($\times 20.66$)
Exp- $k=2$ (7.5, 2.5)	24.78 ($\times 2.12$)	74.69 ($\times 1.92$)	<i>63.86</i> ($\times 1.64$)	82.70 ($\times 1.78$)	<i>76.16</i> ($\times 1.46$)	36.92 ($\times 2.60$)	9.63 ($\times 1.92$)	26.48 ($\times 2.91$)	54.83 ($\times 24.76$)	47.73 ($\times 29.69$)	49.58 ($\times 7.07$)
Exp- $k=2$ (7.5, 0)	24.78 ($\times 2.71$)	74.69 ($\times 1.94$)	<i>63.86</i> ($\times 1.64$)	82.59 ($\times 2.09$)	<i>76.16</i> ($\times 1.67$)	35.94 ($\times 3.22$)	9.30 ($\times 2.56$)	25.81 ($\times 3.73$)	52.43 ($\times 25.23$)	46.36 ($\times 29.89$)	49.19 ($\times 7.52$)
Exp- $k=4$ (7.5, 2.5)	24.78 ($\times 2.68$)	74.69 ($\times 1.93$)	<i>63.86</i> ($\times 1.64$)	82.75 ($\times 2.01$)	<i>76.16</i> ($\times 1.49$)	36.01 ($\times 3.29$)	9.38 ($\times 2.49$)	25.70 ($\times 3.88$)	51.04 ($\times 25.69$)	45.61 ($\times 30.10$)	49.30 ($\times 7.56$)
Exp- $k=4$ (7.5, 0)	24.55 ($\times 3.92$)	75.24 ($\times 2.50$)	<i>63.86</i> ($\times 1.64$)	80.47 ($\times 2.51$)	<i>74.66</i> ($\times 2.54$)	33.74 ($\times 4.57$)	9.29 ($\times 3.80$)	24.03 ($\times 5.59$)	45.49 ($\times 27.83$)	41.09 ($\times 31.67$)	47.83 ($\times 8.77$)
Exp- $k=8$ (7.5, 2.5)	25.22 ($\times 3.44$)	74.66 ($\times 2.07$)	<i>63.86</i> ($\times 1.64$)	82.15 ($\times 2.37$)	<i>76.16</i> ($\times 1.72$)	34.42 ($\times 4.33$)	9.29 ($\times 3.31$)	24.29 ($\times 5.50$)	44.02 ($\times 28.83$)	39.90 ($\times 32.36$)	47.79 ($\times 8.47$)
Exp- $k=8$ (7.5, 0)	25.22 ($\times 5.91$)	74.65 ($\times 4.21$)	<i>63.86</i> ($\times 1.64$)	79.27 ($\times 4.76$)	<i>74.03</i> ($\times 4.64$)	30.33 ($\times 7.03$)	9.14 ($\times 5.74$)	20.86 ($\times 9.53$)	38.49 ($\times 33.37$)	34.79 ($\times 36.03$)	45.46 ($\times 11.19$)
Linear (7.5, 2.5)	25.22 ($\times 1.88$)	74.69 ($\times 1.92$)	<i>63.86</i> ($\times 1.64$)	82.70 ($\times 1.78$)	<i>76.16</i> ($\times 1.45$)	37.38 ($\times 2.31$)	9.99 ($\times 1.74$)	26.73 ($\times 2.42$)	55.36 ($\times 24.66$)	47.92 ($\times 29.65$)	49.88 ($\times 6.85$)
Linear (7.5, 0)	24.78 ($\times 2.27$)	74.69 ($\times 1.92$)	<i>63.86</i> ($\times 1.64$)	82.70 ($\times 1.78$)	<i>76.16</i> ($\times 1.49$)	36.47 ($\times 2.67$)	9.27 ($\times 2.10$)	26.64 ($\times 3.02$)	54.90 ($\times 24.73$)	47.74 ($\times 29.68$)	49.58 ($\times 7.10$)

Table 8: **LLaDA Instruct** with intermediate-curvature exponentials ($k \in \{4, 8\}$) and explicit thresholds. Highest accuracy per column in **bold**; second-highest in *italics*. The rightmost column reports mean accuracy with mean speedup in parentheses.

LLaDA. Results on *LLaDA* align closely with the *Dream* trends for both base and instruction-tuned variants. For *LLaDA Base* (Table 7), Exp- $k=4$ (7.5, 2.5) ties for the highest overall average (53.17) while achieving a $\times 2.13$ speedup; moving to Exp- $k=8$ (7.5, 2.5) increases the speed to $\times 2.55$ with a small average dip (52.55), and relaxing to (7.5, 0) further boosts speed (e.g., Winogrande 78.69 at $\times 3.99$) at the cost of broader quality drops—mirroring the *Dream* trade-off. For *LLaDA Instruct* (Table 8), the curvature-speed relationship is monotone while quality degrades gradually with curvature: averages move from 49.58 at $\times 7.07$ for Exp- $k=2$ (7.5, 2.5) to 49.30 at $\times 7.56$ for $k=4$, 47.79 at $\times 8.47$ for $k=8$, and 47.48 at $\times 10.87$ for $k=16$. In short, the intermediate-curvature exponentials again occupy the expected middle ground between the gradual-decay and high-curvature regimes, providing a convenient knob to trade a bit more speed for a small, predictable loss in fidelity.

D Qualitative Assessment of *SchED*

To complement the aggregate benchmark results, we present representative per-instance examples comparing vanilla diffusion decoding (fixed horizon) to *SchED* under an exponential schedule ($k=2$). Table 9 reports the realized number of denoising steps and wall-clock latency under identical inference settings, alongside the raw model outputs.

Question	Vanilla decoding answer (raw)	<i>SchED</i> answer (raw)	Base steps	<i>SchED</i> steps	Base time (s)	<i>SchED</i> time (s)
What is the capital of France?	The capital of France is Paris. It is a global city known for its cultural heritage and serves as the seat of the French executive, government, as well as other institutions.	The capital of France is Paris.	128	3	5.384	0.130
Translate: "Hello, my name is John."	Bonjour, moi, mon nom est John.	Bonjour., mon nom est John..	128	3	5.533	0.134
Write a python program to calculate the Fibonacci of a given number.	Here's a simple Python program to calculate the Fibonacci of a given number: <pre>def fibonacci(n): if n <= 0: seq = [] elif n == 1: seq = [0] else: seq = [1,1] i = 2 12 12 while i<n+1: seq += seq[-1]+seq[-2] i+=1 return seq 12 print(fibonacci(12)) 12</pre>	def fibonacci(n): if n <= 1: return n else: return fibonacci(n-1) + fibonacci(n-2)	256	237	15.284	14.380
A bookstore had 360 books in stock. In the morning it sold 3/8 of them. How many books are left at the end of the day? Think step by step.	3 / 8 of = 3 / 8 * 360 =135. So, the bookstore has 360 - 135 = 135 books at the end of the day. ### 135. The answer is, 360.	3 sold 3 0 * = 1355.. There left 360 - 135 = 225 books. The answer: 225.	256	235	17.071	16.161
Write a Python program to find the greatest common divisor (GCD) of two numbers.	def gcd(x, y): while(y): x, y=x % y, x num1 = 0 num2 = 0 try: num1 = int(input("Enter the first number and num2:")) print("GCD of: ", num1, " and: ", num2, "are: ", gcd(num1, num2)) except: break	def gcd(a, b): while b: a, b = b, a % b return a	256	238	17.073	16.176
What is the square root of 144?	$\sqrt{144} = 12$	The square root of 144 is $\boxed{12}$	256	128	15.586	8.006

Table 9: Qualitative comparison between vanilla diffusion decoding (full horizon) and *SchED* (exponential schedule, $k=2$). Base steps correspond to the fixed diffusion horizon T_{\max} , while *SchED* steps report the realized number of denoising steps before early commitment. Beyond large speedups on easy prompts, several examples demonstrate *strict correctness improvements*: vanilla decoding can exhibit late-step degradation (incorrect arithmetic or malformed code), while *SchED* preserves a correct intermediate solution at reduced compute.

The examples from Table 9 illustrate two recurring operating regimes consistent with *SchED*'s progress-aware confidence criterion. First, for short and low-complexity prompts, the aggregated logit-margin confidence crosses the schedule threshold within the first few refinement steps (e.g., 3 steps out of 128 or 256), yielding large latency reductions. In this early-exit regime, minor surface artifacts (e.g., extra punctuation) may occasionally appear, while semantic content is typically preserved; such artifacts can be reduced by adopting more constrained schedules (e.g., higher τ_{low} or slower decay). Second, for prompts that require multi-step reasoning or structured generation (e.g., arithmetic word problems and code), *SchED* often runs close to the full horizon (e.g., 235–238 steps out of 256), reflecting sustained uncertainty and avoiding premature termination. Beyond accelerating easy instances, Table 9 also includes cases where *SchED* yields *more correct* outputs than full-horizon diffusion decoding. In these examples, additional refinement steps in vanilla decoding introduce late-stage degradation—e.g., numerical slips, contradictions, or syntactic corruption—despite using the full budget. By terminating once confidence has stabilized over the *full answer span*, *SchED* can preserve a higher-quality intermediate solution. For instance, *SchED* produces the correct arithmetic result for the bookstore problem (225 books remaining) where vanilla decoding outputs an incorrect remainder. Similarly, on code-generation prompts, *SchED* yields clean, executable code-implementations, whereas vanilla decoding degenerates into malformed programs. These examples illustrate that *SchED* is not solely a speed–quality trade-off knob: on some inputs it provides a *Pareto improvement*, reducing compute while improving correctness by avoiding late-step drift.