
Learning Unmasking Policies for Diffusion Language Models

Anonymous Authors¹

Abstract

Diffusion (Large) Language Models (dLLMs) now match the downstream performance of their autoregressive counterparts on many tasks, while holding the promise of being more efficient during inference. One critical design aspect of dLLMs is the *sampling procedure* that selects which tokens to unmask at each diffusion step. Indeed, recent work has found that heuristic strategies such as confidence thresholding improve both sample quality and token throughput compared to random unmasking. However, such heuristics have downsides: they require manual tuning, and we observe that their performance degrades with larger block sizes. In this work, we instead propose to train sampling procedures using reinforcement learning. Specifically, we formalize masked diffusion sampling as a Markov decision process in which the dLLM serves as the environment, and propose a lightweight policy based on a single-layer transformer that maps dLLM token confidences to unmasking decisions. Our experiments show that these trained policies match the performance of state-of-the-art heuristics when combined with semi-autoregressive (block) generation, while outperforming them in the full-diffusion setting.

1. Introduction

Discrete diffusion (Austin et al., 2021a; Hoogeboom et al., 2021; Dieleman et al., 2022; Lou et al., 2024; Shi et al., 2024) has recently emerged as a compelling alternative to the predominant autoregressive (AR) modeling paradigm in large language models (LLMs). Unlike AR models, which generate the next token in a left-to-right fashion (Radford et al., 2018), diffusion LLMs (dLLMs) generate text by learning to reverse a noising process that progressively corrupts token sequences. In particular, masked diffusion

models (MDMs; Sahoo et al. 2024; Ou et al. 2025), a subclass of discrete diffusion models that use time-dependent BERT-style (Devlin et al., 2019) masking as the forward noising process, have recently demonstrated impressive performance, with models like LLaDA (Nie et al., 2025) and Dream (Ye et al., 2025) matching the performance of similarly sized autoregressive LLMs.

At generation time, MDMs begin with a fully masked sequence and iteratively unmask a fixed number of tokens at randomly sampled positions in each sampling step. As a result, they offer the potential for faster inference, as they can, in principle, generate multiple tokens in parallel using a single model call. Despite this promise, open-source dLLMs have, until recently, lagged behind their AR counterparts in terms of inference efficiency. This changed with Fast-dLLM (Wu et al., 2025), which demonstrated that a dLLM like LLaDA (Nie et al., 2025) can achieve higher token throughput than similarly sized LLaMA models (Dubey et al., 2024) while maintaining competitive performance. A key component of Fast-dLLM lies in its sampling heuristic: instead of unmasking a fixed number of randomly sampled token positions, it proposes to unmask all tokens whose confidences exceed a pre-specified threshold. This success has since inspired the development of increasingly sophisticated sampling heuristics (see Appendix A for a comprehensive overview), further advancing the state of the art. These heuristics can, however, be difficult to tune, and seem to work best when imposing semi-autoregressive (semi-AR) generation, in which small blocks of tokens are unmasked sequentially (Arriola et al., 2025a; Nie et al., 2025). We posit that such limitations are difficult to avoid with heuristics alone, since they are essentially handcrafted solutions to a sequential decision making problem: *What is the optimal order in which to unmask the sequence of tokens, in order to strike a good balance between correctness and efficiency?*

We address this question by moving beyond heuristics, instead proposing a learning-based approach that leverages a transformer-based unmasking policy trained via reinforcement learning (RL). This approach is motivated by the above observation that unmasking in dLLMs can be viewed as a (Markovian) sequential decision making problem, and follows recent successes in RL for post-training of language models. However, in contrast to recent works that use RL to improve the reasoning abilities of dLLMs (e.g., Zhao et al.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

2025b), our goal is not to improve the capacity of the underlying model, but rather to use RL to automate the discovery of adaptive sampling strategies. We therefore treat the underlying dLLM as the *environment* in which to act, not the policy, leaving it unchanged and instead parameterizing a policy as a lightweight stand-alone network (Figure 8). Empirically, we demonstrate that our learned sampling policies match the performance of state-of-the-art heuristic samplers (Wu et al., 2025) in standard semi-AR generation, while also surpassing heuristics in full-diffusion setting. In summary, we make the following contributions:

- We formalize sampling in dLLMs as a Markov decision process (MDP) (Section 3.1), propose a lightweight design for the unmasking policy (Section 3.2), and outline our RL training pipeline based on group relative policy optimization (GRPO; Shao et al. 2024) (Section 3.3).
- Our experiments show that learned policies with our RL framework match the performance of heuristic samplers such as Fast-dLLM (Section 4.1), while also addressing some of the challenges heuristic samplers face outside the semi-AR regime (Section 4.2).
- We study the transferability of our learned samplers across different models and data domains (Section 4.3), and provide detailed insights into stabilizing RL training and policy design (Section 4.4). We also demonstrate that policies exhibit qualitatively different unmasking behavior compared to confidence heuristics (Appendix D).

2. Background

Throughout the paper, we use the notation $[L]$ to represent the set $\{1, \dots, L\}$. We denote a sequence of tokens as $\mathbf{x} = (x^1, \dots, x^d) \in \mathcal{V}^d$, where d is the sequence length and $\mathcal{V} := [V]$ is the vocabulary. Given our focus on MDMs, we assume the vocabulary includes a special mask token M .

2.1. Masked Diffusion Models

We focus on masked diffusion models (MDMs), deferring a more general introduction to discrete diffusion to Appendix H.

Training. MDMs learn to generate data by training a BERT-style (Devlin et al., 2019) masked predictor to reverse the forward noising process. Concretely, given a training sample $\mathbf{x}_0 \sim p_{\text{data}}$, the forward process corrupts each token independently by setting it to the mask token with probability proportional to the diffusion timestep $t \in [0, 1]$:

$$p_t(\mathbf{x}_t | \mathbf{x}_0) := \prod_{k=1}^d p_t(x_t^k | \mathbf{x}_0),$$

$$\text{where } p_t(x_t^k | \mathbf{x}_0) := \begin{cases} 1 - t, & \text{if } x_t^k = x_0^k, \\ t, & \text{if } x_t^k = M, \\ 0, & \text{otherwise.} \end{cases}$$

Recent work has shown that an MDM parametrized by θ can learn the reverse process by maximizing the following evidence lower bound (ELBO) (Ou et al., 2025), $\mathcal{L}(\theta) \leq \mathbb{E}_{\mathbf{x}_0 \sim p_{\text{data}}}[\log p_{\theta}(\mathbf{x}_0)]$, denoting $\mathbf{1}[\cdot]$ the indicator function:

$$\mathcal{L}(\theta) := \mathbb{E}_{t \sim U[0,1], \mathbf{x}_0 \sim p_{\text{data}}, \mathbf{x}_t \sim p_t(\mathbf{x}_t | \mathbf{x}_0)} \left[\frac{1}{t} \sum_{k=1}^d \mathbf{1}[x_t^k = M] \log p_{\theta}(x_0^k | \mathbf{x}_t) \right]. \quad (1)$$

Generation. When generating an answer for a given prompt \mathbf{x} , MDMs start with a sequence of all-masked tokens $\mathbf{y}_T := (M, \dots, M)$, where $L := |\mathbf{y}_T|$ is a pre-specified maximum answer length. Then, in each sampling step $t \in [T]$, the MDM outputs token distributions $p_t^k := p_{\theta}(y_0^k = \cdot | \mathbf{x}, \mathbf{y}_t)$ for all token positions $k \in [L]$ and a sampling strategy decides which subset $\mathcal{U}_t \subseteq \mathcal{M}_t$ of the still-masked positions $\mathcal{M}_t := \{k \in [L] | y_t^k = M\}$ to unmask. A new, partially unmasked/denoised answer \mathbf{y}_{t-1} is obtained via

$$\mathbf{y}_{t-1} := \begin{cases} y \sim p_t^k, & \text{if } k \in \mathcal{U}_t, \\ y_t^k, & \text{otherwise.} \end{cases} \quad (2)$$

The stochasticity of sampling $y \sim p_t^k$ depends on the dLLM temperature τ , with higher temperatures leading to more diverse (but often lower quality) generations. When evaluating dLLMs, it is common to set $\tau = 0$ (Nie et al., 2025; Wu et al., 2025), resembling greedy decoding. The choice of sampling/unmasking strategy is also important, since different choices will produce different unmasking sets \mathcal{U}_t , thus affecting both sampling quality and efficiency. This has drawn considerable attention to the development of optimal sampling techniques for dLLMs (see Appendix A).

2.2. Heuristic Samplers

One popular approach to decide which positions \mathcal{U}_t to unmask at each timestep is to construct a heuristic that leverages uncertainty measures derived from the predicted token distributions p_t^k (Wu et al. 2025; Ben-Hamu et al. 2025; *inter alia*). In particular, recent work has shown substantial efficiency gains with heuristic strategies that employ the *confidence* $c_t^k := \max_{v \in \mathcal{V}} p_t^k(v)$ of the underlying dLLM in order to make their decisions. Two representative methods within this space are *high-confidence* unmasking (Chang et al., 2022), in which each forward pass unmask a fixed number of tokens (K) with the highest confidences,

$$\mathcal{U}_t^K := \left\{ \arg \max_{I \subseteq \mathcal{M}_t, |I|=K} \sum_{k \in I} c_t^k \right\}$$

as well as the confidence-thresholding strategy of Fast-dLLM (Wu et al., 2025), which allows for a variable number of tokens to be unmasked at each step by comparing the confidences to a fixed threshold λ :

$$\mathcal{U}_t^{\lambda} := \{k \in \mathcal{M}_t | c_t^k > \lambda\}.$$

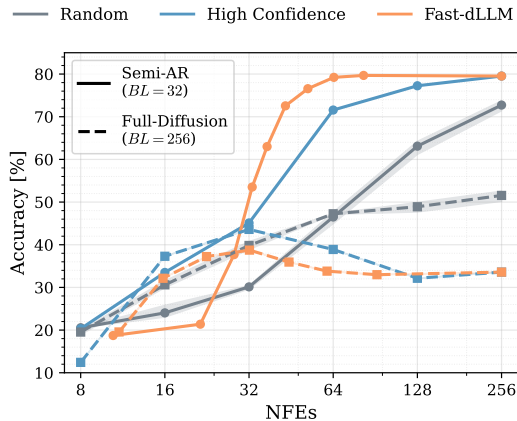


Figure 1. LLaDA-8B-Instruct (Nie et al., 2025) on GSM8k, with semi-AR generation ($BL = 32$; —) and without (full-diffusion regime, $BL = L = 256$; - -). More datasets and models in Appendix C.1. Generation speed is measured in network function evaluations (NFEs), which corresponds to the number of sampling steps.

Despite the undeniable success of confidence-based samplers, relying on handcrafted solutions comes with drawbacks. Beyond the obvious need for careful design (e.g., selecting an appropriate confidence measure or threshold λ), we also found that these methods often exhibit high sensitivity to the exact sampling configuration. For instance, many heuristics rely on semi-autoregressive (semi-AR) generation (Arriola et al., 2025a; Nie et al., 2025), where diffusion decoding proceeds one “block” of contiguous tokens at a time, with block length BL serving as yet another hyperparameter. As shown in Figure 1, outside this semi-AR regime, the performance of confidence-based heuristics can degrade below that of a random unmasking, and no longer benefit from increased compute (i.e., higher numbers of function evaluations, NFEs). These limitations raise the question of whether effective sampling methods can be *learned* rather than manually tuned.

3. Learning Unmasking Policies

We now introduce our approach to learn sampling in dLLMs via reinforcement learning (RL). We begin by formulating sampling as a Markov decision process (MDP) (Section 3.1), then propose a sampling policy (Section 3.2), and finally provide details on our RL training (Section 3.3).

3.1. dLLM Sampling as a Markov Decision Process

To facilitate the development of RL-based samplers, we describe first the Markov decision process (MDP) (Sutton et al., 1998) for sampling in dLLMs. To stay aligned with the MDP notation in Section 2.1, we reverse the time in our MDP formulation: $t = T, T - 1, \dots, 1$, with T denoting the maximum horizon (number of sampling steps).

- The *state* $(\mathbf{x}, \mathbf{y}_t)$ holds a prompt $\mathbf{x} \in \mathcal{V}^d$ and the current (partially masked) generation $\mathbf{y}_t \in \mathcal{V}^L$. For brevity, we omit \mathbf{x} from the state notation unless needed. The initial state has a fully masked generation: $\mathbf{y}_T = (\mathbf{M}, \dots, \mathbf{M})$.
- An *action* $\mathbf{u}_t \in \{0, 1\}^L$ is a vector of unmasking decisions indicating which positions have been selected to be unmasked in the next transition step. These actions are chosen according to the policy π_ϕ : $\mathbf{u}_t \sim \pi_\phi(\cdot | \mathbf{y}_t)$ (introduced later in Section 3.2).
- The *transition* $\mathbf{y}_{t-1} \sim P(\mathbf{y}_{t-1} | \mathbf{y}_t, \mathbf{u}_t; \tau)$ corresponds to standard sampling in dLLMs (cf. Equation (2)), with the action \mathbf{u}_t determining which tokens get unmasked: $\mathcal{U}_t^\pi := \{k \in \mathcal{M}_t | u_t^k = 1\}$.
- The *reward* $R(\mathbf{y}, \mathbf{y}_t)$ is provided only at the final generation step, which corresponds to the first timestep where all tokens have been unmasked $\hat{T} := \max\{t \in [T] | y_t^k \neq \mathbf{M}, \forall k \in [L]\}$. To learn useful samplers, the reward should promote both correctness (i.e., the generated answer $\mathbf{y}_{\hat{T}}$ being ‘close’ to the reference answer \mathbf{y}) and efficiency (i.e., minimizing the number of steps $T - \hat{T}$). Note that the reward depends on action \mathbf{u}_t , and therefore on the policy π_ϕ , implicitly through its influence on the generations \mathbf{y}_t and the number of steps. We defer our concrete choice of the reward function to Section 3.3.

With the MDP above, finding the optimal sampling becomes a standard RL problem of finding a policy π_ϕ parametrized by ϕ that maximizes the expected reward:

$$\max_{\phi} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p_{\text{data}}, \mathbf{u}_t \sim \pi_\phi, \mathbf{y}_t \sim P} \left[\sum_{t=1}^T R(\mathbf{y}, \mathbf{y}_t) \right].$$

3.2. Lightweight Confidence Policy Design

Having outlined the MDP above, we next describe our proposed implementation for the sampling policy π_ϕ (see Appendix B for policy sampling diagram and algorithm).

Confidence-based input. Recall that a state consists of a partially masked token sequence \mathbf{y}_t . To avoid constructing policies that operate at the token level, and thereby minimize computational overhead, we rely on the vector of token confidences $\mathbf{c}_t := (c_t^1, \dots, c_t^L)$, which are readily available since the token-level predictive distributions p_t^k are computed at each transition step $\mathbf{y}_{t-1} \sim P$ anyways. Our choice is further motivated by the aforementioned heuristic methods that primarily operate on confidences (Wu et al., 2025), and also based on our own ablations (see Section 4.4). Specifically, we observed that alternatively relying on dLLM’s hidden states required significantly larger policy models without consistently improving performance.

Lightweight transformer. We introduce a small, learnable neural network f_ϕ that maps the vector of confidences \mathbf{c}_t to a vector of unmasking scores (logits) $\mathbf{b}_t = f_\phi(\mathbf{c}_t, \mathbf{m}_t, t)$ with $\mathbf{b}_t \in \mathbb{R}^L$. We additionally include a binary mask vector $\mathbf{m}_t := (m_t^1, \dots, m_t^L)$, where $m_t^k := \mathbb{1}[k \in \mathcal{M}_t]$ indicates whether position k is still masked, and inform the policy with the time index $t \in [T]$. In practice, f_ϕ is implemented as a lightweight transformer (see Section 4 for more details on the exact policy architecture). Notably, its size is less than 0.01% of the pretrained dLLMs used in our experiments, resulting in negligible computational overhead during sampling (see Figure 10).

Bernoulli likelihood. We sample the unmasking actions via $u_t^k \sim \text{Ber}(s_t^k)$ with $s_t^k := \sigma(b_t^k)$, where $\text{Ber}(\cdot)$ denotes the Bernoulli distribution and $\sigma(\cdot)$ is the sigmoid function. Conveniently, the policy likelihood $\pi_\phi(\mathbf{u}_t) := \prod_{k=1}^L (s_t^k)^{u_t^k} \cdot (1 - s_t^k)^{1-u_t^k}$ is readily available in closed form and does not require any additional approximations (unlike in the case of post-training dLLMs; Zhao et al. 2025b). We also considered an alternative formulation based on the Plackett-Luce model (Luce, 1959; Plackett, 1975) which we call *dynamic Plackett-Luce sampling* (DPLS; see Appendix F), but since our ablations showed comparable performance (Section 4.4), we favor the Bernoulli formulation in our experiments for its simplicity.

At generation, we additionally propose to “temper” the Bernoulli probabilities as $s_t^k := \sigma(b_t^k / \tau_\pi)$. The temperature τ_π thus controls the sharpness of the policy distribution; we found that this can sometimes serve as a useful test-time knob for trading off accuracy versus efficiency by making the policy more or less “decisive” (Section 4).

To ensure convergence when all sampled actions are zero ($\mathbf{u}_t = \mathbf{0}$), we unmask the position with the highest Bernoulli parameter s_t^k (similar to Fast-dLLM). We apply this fallback only at test time, as we found that forcing unmasking during training was prone to reward hacking. Instead, if not all tokens are unmasked within T steps during training, we simply terminate generation and evaluate the answer with some tokens left unmasked. Since such samples tend to get low or no reward, we find that the policies quickly learn not to leave any tokens unmasked without explicit instruction.

3.3. Learning π_ϕ with GRPO

To train our sampling policy, we adopt group relative policy optimization (GRPO) (Shao et al., 2024), a method recently popularized as a simpler and more scalable alternative to earlier policy gradient approaches such as PPO (Schulman et al., 2017). Specifically, for each prompt $\mathbf{x} \in \mathcal{D}$, we sample G trajectories of generations $\{\mathbf{y}_T^g, \dots, \mathbf{y}_{\hat{T}_g}^g\}_{g=1}^G$ along with their corresponding unmasking decisions $\{\mathbf{u}_T^g, \dots, \mathbf{u}_{\hat{T}_g}^g\}_{g=1}^G$. Importantly, we fix the

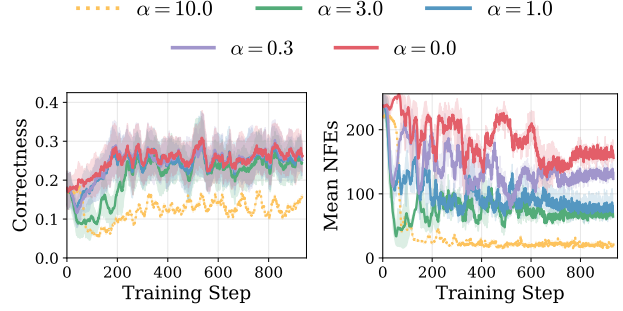


Figure 2. Correctness reward (rolling average, 20 steps) on GSM8k (left) and average number of sampling steps (right) during training of our policies for various values of α (cf. Equation (3)). Averaged over two random seeds, with shaded areas indicating (min, max); only one seed shown for $\alpha = 10.0$ (\dots) due to training instability.

dLLM sampling temperature τ to 0 (i.e., greedy decoding) to ensure that any variation among samples within a group arises solely from different unmasking actions. After computing rewards for each sample in the group, we define the advantage as $A_t^g := R(\mathbf{y}, \mathbf{y}_{\hat{T}_g}^g) - \frac{1}{G} \sum_{i=1}^G R(\mathbf{y}, \mathbf{y}_{\hat{T}_g}^i)$, following recent best practices that recommend omitting standard deviation normalization of advantages (Zhao et al., 2025b). Additionally, note that the reward at the final generation step \hat{T}_g for each sample is propagated to all preceding timesteps t to provide a learning signal throughout the entire sampling process. Our final training objective is then

$$\mathcal{J}(\phi) := \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}, \{\mathbf{y}_{T:\hat{T}_g}^g\}_{g=1}^G \sim P_\theta, \{\mathbf{u}_{T:\hat{T}_g}^g\}_{g=1}^G \sim \pi_{\phi_{\text{old}}}} \left[\frac{1}{G} \sum_{g=1}^G \frac{1}{T - \hat{T}_g} \sum_{t=\hat{T}_g}^T \min \left\{ \rho_t^g \cdot A_t^g, \text{clip}(\rho_t^g, 1 - \epsilon, 1 + \epsilon) \cdot A_t^g \right\} \right]$$

where π_ϕ is the current policy being updated, and $\pi_{\phi_{\text{old}}}$ refers to an earlier version of the policy used to generate the RL rollouts. The likelihood ratio $\rho_t^g := \frac{\pi_\phi(\mathbf{u}_t^g)}{\pi_{\phi_{\text{old}}}(\mathbf{u}_t^g)}$ serves as the importance sampling correction term, which together with clipping (via ϵ) aims to stabilize the off-policy training. Note that we exclude already-unmasked positions from the policy likelihood computation. Finally, since our policies are trained from scratch, we remove the KL regularization term in our GRPO objective.

Reward As mentioned in Section 3.1, we wish to obtain a policy that yields ‘correct’ generations while also being fast. To this end, we define our final *multiplicative* reward as

$$R(\mathbf{y}, \mathbf{y}_t) := \begin{cases} r(\mathbf{y}, \mathbf{y}_t) \cdot \left(1 - \frac{T-t}{T}\right)^\alpha, & \text{if } t = \hat{T}, \\ 0, & \text{otherwise,} \end{cases} \quad (3)$$

where $r(\mathbf{y}, \mathbf{y}_t)$ is a task-specific correctness term (e.g., a binary reward indicating whether the generated mathemat-

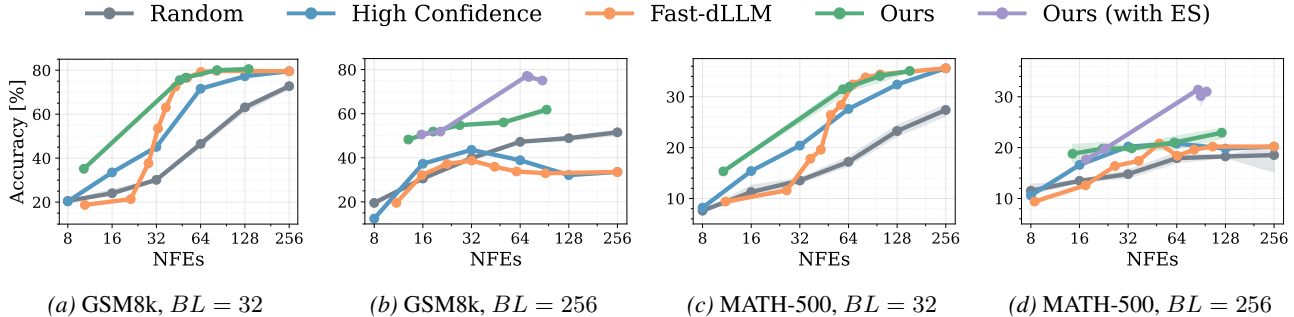


Figure 3. Results for LLaDA in semi-AR (Figure 3a & Figure 3c) and full-diffusion (Figure 3b & Figure 3d) generation regimes ($L = 256$). Results for Dream-7B are provided in Figure 11. For our policies we vary $\alpha \in \{10, 3, 1, 0.3, 0\}$ and use $\tau_\pi = 0.5$ for $BL = 32$ and $\tau_\pi = 1$ for $BL = 256$. Expert steering (ES) described in detail in Appendix G. Wall-clock time plots shown in Figure 10.

ical answer is correct). To encourage faster sampling, we incorporate a computational penalty based on the number of steps, $T - \hat{T}$, with $\alpha \geq 0$ serving as a hyperparameter that controls the trade-off between accuracy and speed. While we experimented with an additive penalty of the form $r(\mathbf{y}, \mathbf{y}_{\hat{T}}) - \alpha \left(\frac{T - \hat{T}}{T} \right)$ (Graves, 2016) we found that this led to problematic reward hacking: when all samples in a group are incorrect, as is the case early on when training from scratch, faster samples may still receive a positive advantage, despite being wrong (cf. Section 4.4).

4. Experiments

We begin by verifying that the learned policies match the performance of confidence-based heuristics (Section 4.1). We then highlight the potential of RL sampling to outperform heuristics in the full-diffusion setting (Section 4.2), and examine the transferability of RL policies across models, datasets, and generation lengths (Section 4.3). Next, we explore different ways to instantiate the MDP and analyze their effects on performance (Section 4.4). Finally, we study the unmasking behavior of RL policies, highlighting their qualitative differences to heuristics (details in Appendix D).

4.1. Learning Effective Sampling via RL

We start by demonstrating that our proposed RL framework yields effective sampling strategies in dLLMs, where effectiveness is defined in terms of both accuracy and speed.

Experiment details. We use LLaDA-8B-Instruct (Nie et al., 2025) as the base dLLM (additional results on Dream-7B-Instruct (Ye et al., 2025) provided in Appendix C.3). We parametrize the policy network f_ϕ as a shallow (single-layer) transformer incorporating adaptive layer normalization for conditioning; hyperparameter details for the architecture are provided in Appendix I and an architecture diagram can be found in Appendix J. We train five different policies, corresponding to $\alpha \in \{10, 3, 1, 0.3, 0\}$. Each policy is trained semi-autoregressively at $BL = 32$ on a

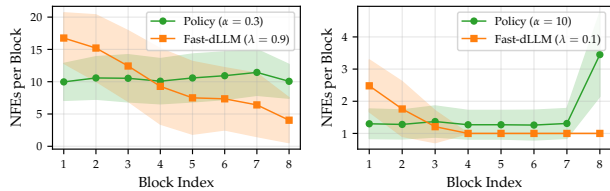
single epoch of mixture data, sampled proportionally from the training sets of GSM8k (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021), resulting in roughly 15,000 training samples. Figure 2 shows the training dynamics; as expected, higher α generally gives rise to faster policies.

We then compare the resulting policies on the test sets of both GSM8k and MATH to the confidence-based heuristics introduced in Section 2.2, picking one out of two training seeds based on the final training loss and averaging over three test-time seeds. To obtain Pareto frontiers for the heuristic methods, we use $K \in \{8, 16, 32, 64, 128, 256\}$ for the random baseline and high-confidence unmasking, while for Fast-dLLM we use $\lambda \in \{0.1, 0.2, \dots, 1.0\}$ throughout. For all methods, we use the standard greedy decoding setting ($\tau = 0$) when generating test answers.

Results for (Short) $BL = 32$, Figure 3a and Figure 3c.

Generally, we observe that the performance of the learned policies (\rightarrow) exceeds those of both the random baseline (\rightarrow) as well as the high-confidence sampling (\rightarrow), while matching Fast-dLLM (\rightarrow). That our learned policies do not surpass Fast-dLLM in the mid-to-high NFEs range may suggest that this heuristic is near-optimal under semi-AR generation. Interestingly, despite very similar quantitative performance, we find that the learned policies exhibit qualitatively somewhat different unmasking behavior compared to Fast-dLLM (see Appendix D). Concretely, the policy un-masks adjacent tokens much less frequently (Figure 20) and allocates compute more uniformly across blocks (Figure 4a).

We note that the effect of α in training (cf. Figure 2) carries over to the test set: policies trained with higher values of α exhibit faster, yet less accurate, sampling. Of particular note is $\alpha = 10.0$, which exhibited greater training instability due to the very sharp slope of the reward. When successfully trained, this yields a very fast policy that outperforms Fast-dLLM in the low-NFE (~ 10) regime on both datasets. Qualitatively, we find evidence that this improved performance is likely due to the policy learning to slow down in the final block when generating a numerical answer, as



(a) BL32-slow (○)

(b) BL32-fast (○)

Figure 4. Average number of sampling steps (NFEs) per block for LLaDA on GSM8k with semi-AR generation (same setting as in Figure 3a). *Left* we show NFEs per block for ‘slow’ variants of Fast-dLLM and policy sampling (~ 75 NFEs), while in the *right* plot we show ‘fast’ variants (~ 10 NFEs). Fast-dLLM exhibits a pattern of allocating more compute to earlier blocks. Our policy sampling, by contrast, distributes compute more uniformly across blocks, except in the fast ($\alpha = 10$) policy, where most compute is expended in the final block while generating numerical answers. To see where evaluated policies here are located on the pareto frontier, see Figure 17a. More details on qualitative differences between Fast-dLLM and policy sampling are provided in Appendix D.

shown in Figure 4b (see also Appendix D for more details). This highlights the potential of RL-based policies when optimizing for maximal efficiency.

However, RL policies appear to exhibit less *controllability*, as varying α results in a less smooth traversal of the Pareto frontier than varying the confidence threshold λ . For example, the behavior of the $\alpha = 3.0$ policy is nearly identical to that of $\alpha = 1.0$, despite the reward incorporating a computational penalty that scales exponentially with α . Furthermore, we find that this cannot easily be remedied by using a tighter α grid. In Appendix C.4 we let α range in $\{10.0, 9.0, \dots, 1.0, 0.3, 0.0\}$, but observe that using $\alpha \geq 4.0$ consistently results in converging to a policy either roughly equivalent to that of $\alpha = 3.0$ or to that of $\alpha = 10.0$, with no value of α successfully yielding a policy which interpolates between the two extremes. Besides α , we observe that one can make small changes to the trade-off between compute and performance by varying the policy temperature τ_π at test time; we detail the impact of this parameter in Figure 14.

4.2. Beyond Semi-Autoregressive Decoding

As mentioned in Section 2.2, heuristic samplers often rely on semi-AR generation to achieve good performance. While not typically thought of as problematic in textual domains where strong AR dependencies exist, semi-AR approaches can only partially fulfill the promise of fully parallel generation by restricting the unmasking to the current block. Hence, we train here policies without semi-AR generation, targeting the full-diffusion task ($BL = L = 256$) at both training and test time. We keep all other implementation details identical to those in the previous section.

Results for (Long) $BL = 256$, Figure 3b & Figure 3d. Although all sampling strategies experience a performance

drop relative to the semi-AR setting (cf. Figure 3a & Figure 3c), we find that the policies produced by RL (—) exhibit the smallest decline and consequently achieve the best overall performance. Furthermore, the results suggest that this methodology is able to produce policies which achieve solid performance even in the low-NFE regime. In particular, the learned policies obtain $\sim 50\%$ accuracy at ~ 12 NFEs on GSM8k (compared to $\leq 30\%$ for the heuristic methods regardless of semi-AR use), highlighting the potential of non-semi-AR sampling for achieving maximal efficiency gains. Note however that, in theory, an RL policy trained with $BL = L$ could learn to emulate semi-AR sampling on its own. Thus, the fact that policies trained with $BL = L$ underperform those trained with $BL = 32$ in the mid-high NFE range (comparing, for example, Figure 3a to Figure 3b) suggests that the policies remain far from optimal.

One hypothesis for why this is the case is that our training procedure is not encouraging sufficient exploration, instead converging to locally optimal policies. To address this, we encourage further exploration by using samples generated via semi-AR sampling from Fast-dLLM. We refer to this approach as *expert steering*, and describe it in more detail in Appendix G. As shown in Figure 3, with expert steering (—) RL is able to discover policies that almost close the performance gap to the best accuracy achieved in the semi-AR setting at mid-to-high NFEs (e.g., $\sim 80\%$ on GSM8k and $\sim 35\%$ on MATH), while mostly retaining their strong performance in the low-NFE regime. However, we do find that expert steering introduces significant instability during training, further reducing the controllability through α , with multiple values of α collapsing to near-identical policies. We leave further investigations into stabilizing expert steering training for future work. Finally, we observe that policy sampling and Fast-dLLM induce markedly different unmasking orders in the full-diffusion setting (see Figure 5

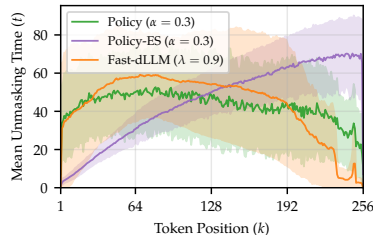
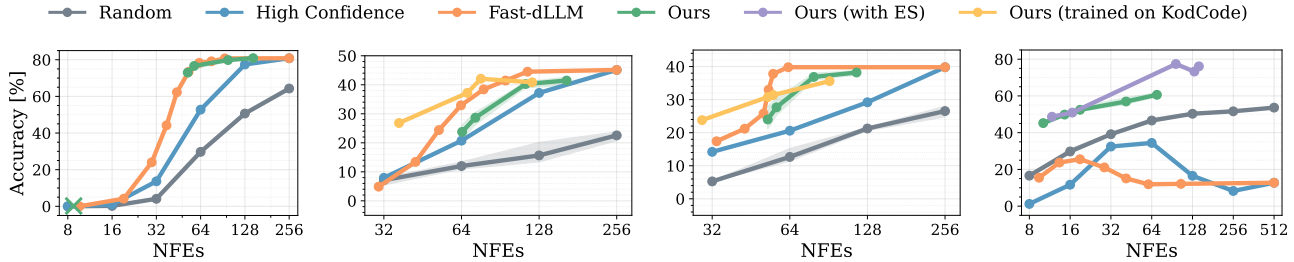


Figure 5. Mean unmasking time for each token position ($L = 256$), averaged over $N = 100$ samples, for LLaDA on GSM8k under full-diffusion generation (same setting as in Figure 3b). For visualization purposes, time is shown in reverse. Fast-dLLM (—) exhibits somewhat reverse (i.e., right-to-left) generation due to LLaDA’s corrupted confidence on padding tokens (see Appendix D). Encouragingly, the expert steering policy (—) learns to overcome this issue and recovers left-to-right generation (observe how tokens at earlier positions are generated first on average). The locations of the evaluated policies on the Pareto frontier are shown in Figure 17b (see ○).



(a) Model transfer: LLaDA \rightarrow Dream on GSM8k. (b) Domain transfer: math \rightarrow coding (HumanEval). (c) Domain transfer: math \rightarrow coding (MBPP). (d) Sequence length transfer: 256 \rightarrow 512 (GSM8k).

Figure 6. Results for the transferability experiments. Note that in (a), the $\alpha = 10$ policy is represented separately by a (x) in the lower left to avoid misleading visualization when interpolating to $\alpha = 3$. For results on coding datasets in (b) and (c), we omit the low-NFE regime, as all approaches degrade to near-zero performance in this setting.

and Appendix D). While Fast-dLLM exhibits reverse (right-to-left) unmasking, the policy with expert steering learns to overcome this behavior and unmask earlier positions first on average (left-to-right), which likely explains its superior performance.

4.3. Transferability of RL Sampling Policies

We next turn to the question of transferability. A notable advantage of heuristic approaches is that they can be applied *post-hoc* to any model or dataset.¹ This naturally raises the question: to what extent can RL policies trained on a specific model and dataset be reused in other settings?

Model transfer: LLaDA \rightarrow Dream. We begin by investigating the transferability of RL policies across models; note that such transfers are possible because our policies rely solely on token confidences and are therefore agnostic to token embeddings, which would not transfer from one model to another. We thus reuse the policies trained on LLaDA (cf. Section 4.1) and evaluate them on Dream; the results are shown in Figure 6a (see Figure 15 for MATH). Encouragingly, most of the LLaDA-trained policies nearly match the performance of Fast-dLLM when evaluated on Dream, and perform very similarly to those trained on Dream directly (Figure 11). The one exception is the $\alpha = 10$ policy (x), which collapses to Fast-dLLM performance and fails to retain the good low-NFE performance observed with LLaDA (cf. Figure 3a)—suggesting that the extreme steepness of the reward curve (due to high α) causes this policy to overfit to model-specific patterns in LLaDA’s confidence levels.

Domain transfer: Mathematical reasoning to code. Next, we examine how well our policies transfer across different domains. We again reuse the policies from Section 4.1, which were trained on a mixture of mathematical data (GSM8k/MATH), but this time evaluate them on the

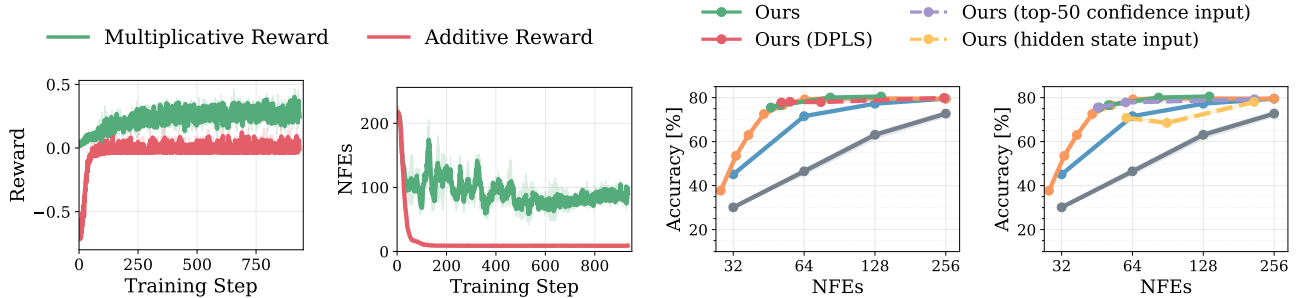
¹Though heuristics might still require manual hyperparameter tuning, as exemplified by different optimal threshold across datasets for Fast-dLLM (e.g., comparing Figure 3a vs Figure 6b).

coding tasks of HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021b). As shown in Figure 6b and Figure 6c (—), we find that these policies fail to fully transfer to the coding domains, especially on HumanEval. To investigate whether this drop in performance is due to the lack of domain-relevant training data, we train a new policy on the coding dataset KodCode-RL-10K (Xu et al., 2025) (—). These coding policies narrow the gap to the Fast-dLLM baseline on HumanEval and improve low-NFE performance on MBPP, underscoring the importance of using a diverse data mixture to support generalization across domains.

Sequence length generalization: $L = 256 \rightarrow L = 512$. Lastly, we examine whether our policies transfer across different sequence lengths L (not to be confused with the block length BL in semi-AR). Such transfer is possible since we instantiate f_θ with a transformer, rather than, for example, an MLP (which requires fixed-length inputs). We thus take the policies from Section 4.2 and evaluate them at a 2x longer sequence length ($L = 512$). Results are shown in Figure 6d (see Figure 16 for MATH). While the baselines degrade further as the sequence length increases, the learned policies yield similar performance to before, suggesting that RL policies can transfer effectively across generation lengths without retraining.

4.4. Exploring the Design Space of π_ϕ

Additive vs. multiplicative rewards. We begin by ablating the structure of the reward function, comparing our proposed multiplicative combination of correctness and computational terms (cf. Equation (3)) with an additive alternative: $r(\mathbf{y}, \mathbf{y}_{\hat{T}}) - \alpha \left(\frac{T - \hat{T}}{T} \right)$. As shown in Figure 7a, while both exhibit increasing reward as training progresses, we observe that the additive formulation is much more prone to ‘reward hacking’, where it collapses to a very-fast-but-often-wrong policy. This is illustrated in Figure 7b: training with the additive reward results in a policy that unmask everything at once, leading the number of sampling steps to collapse to the minimum possible for all inputs. As discussed in



(a) Training reward for LLaDA with additive vs. multiplicative reward function (both $\alpha = 1.0$). (b) Mean NFEs when training LLaDA with additive vs. multiplicative reward (both $\alpha = 1.0$). (c) Bernoulli vs DPLS sampling. Both for LLaDA on GSM8k. (d) Bernoulli policies with varying inputs. All for LLaDA on GSM8k.

Figure 7. Ablations for our proposed RL framework.

Section 3.3, we attribute this issue to the fact that incorrect but fast samples may still receive a positive advantage under the additive reward. In contrast, we find that the multiplicative reward effectively mitigates this issue by assigning a positive advantage only if the generation is correct, resulting in more stable and predictable training behavior.

Policy likelihood. Recall from Section 3.2 that we model the unmasking probability for each position using a Bernoulli distribution. This has the advantage of admitting very efficient inference and likelihood calculations, but relies on the network f_ϕ to implicitly embed relationships between different tokens in the scores \mathbf{b}_t , and runs the risk of producing $\mathcal{U}_t^\pi = \emptyset$ in case $\mathbf{u}_t = \mathbf{0}$. Here, we therefore investigate a more involved sampling procedure, which we call *dynamic Plackett-Luce sampling* (DPLS; detailed description in Appendix F) which is guaranteed to unmask at least one position in each step and which combines the scores \mathbf{b}_t through a softmax. We retrain the policies from Section 4.1 using DPLS; the resulting downstream accuracy on GSM8k is then shown in Figure 7c. We observe that both methods achieve very similar performance, aligning closely with the Fast-dLLM, which might provide additional support for the hypothesis that this frontier is optimal in the semi-AR setting. Furthermore, DPLS policies appear to show slightly better controllability via α (as indicated by a larger spread of policies trained with varying α).

Confidence-based policy input. Finally, we revisit our design choice of relying solely on the maximum token confidence values c_t^k as input to the sampling policy (cf. Section 3.2). We first train and evaluate a set of policies which do not take only the highest confidence per position c_t^k as input, but rather the top 50 highest values, thus giving the model more detailed information about the token predictive distributions p_t^k and potentially allowing it to design its own confidence measures for effective sampling. The results are presented in Figure 7d. Somewhat surprisingly, using only the maximum confidence per position (\rightarrow) performs slightly better than using the top-50 confidences (\leftarrow), which

suggests that alternative uncertainty measures are unlikely to yield performance gains over the simple confidences c_t^k .

Additionally, we explore parametrizing the policy as an additional classification head on top of LLaDA’s final hidden state $\mathbf{h}_t^k \in \mathbb{R}^H$. While this results in a significantly larger policy (300M parameters), it offers the potential advantage of incorporating token-level semantic information. However, in practice we observe that hidden state-based policies perform worse (see \dashrightarrow in Figure 7d) and exhibit less stable training dynamics than the confidence-based policies. These results suggest that the unembedding matrix $W \in \mathbb{R}^{H \times V}$, which maps hidden states to token logits, plays a vital role in enabling effective policy decisions via confidence signals.

5. Conclusion

We introduced a reinforcement learning approach for learning unmasking strategies in diffusion LLMs. Our experiments demonstrated that the learned policies can match or even exceed the performance of recently proposed sampling heuristics, paving the way for the automated discovery of scalable and robust sampling mechanisms.

Limitations and future work. One notable limitation of our approach is that it requires training a separate policy for each α , and that this hyper-parameter sometimes does not yield as much fine-grained control over the policy’s behavior as desired (Section 4.1). Thus, an important future direction is to investigate whether the accuracy-speed trade-off could be controlled through some alternative, ideally test-time, mechanism. We also observe that the optimal policy temperature τ_π varies across generation settings (Figure 14), suggesting the need to explore whether τ_π can be learned jointly with the policy. Beyond overcoming these limitations, promising future extensions are (i) expanding the training mixture to incorporate data from multiple domains; (ii) extending our sampling policies to also support remasking (Wang et al., 2025b; Huang et al., 2025c), and (iii) moving beyond text to learn samplers for multimodal discrete diffusion (Swerdlow et al., 2025; Zhou et al., 2025).

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

Arora, D. and Zanette, A. Training language models to reason efficiently. *NeurIPS*, 2025.

Arriola, M., Gokaslan, A., Chiu, J. T., Yang, Z., Qi, Z., Han, J., Sahoo, S. S., and Kuleshov, V. Block diffusion: Interpolating between autoregressive and diffusion language models. *ICLR*, 2025a.

Arriola, M., Schiff, Y., Phung, H., Gokaslan, A., and Kuleshov, V. Encoder-decoder diffusion language models for efficient training and inference. *NeurIPS*, 2025b.

Austin, J., Johnson, D. D., Ho, J., Tarlow, D., and Van Den Berg, R. Structured denoising diffusion models in discrete state-spaces. *NeurIPS*, 2021a.

Austin, J., Odena, A., Nye, M., Bosma, M., Michalewski, H., Dohan, D., Jiang, E., Cai, C., Terry, M., Le, Q., et al. Program synthesis with large language models. *arXiv:2108.07732*, 2021b.

Azangulov, I., Pandeva, T., Prasad, N., Zazo, J., and Karmalkar, S. Parallel sampling from masked diffusion models via conditional independence testing. *arXiv:2510.21961*, 2025.

Bansal, P. and Sanghavi, S. Enabling approximate joint sampling in diffusion lms. *arXiv:2509.22738*, 2025.

Bao, W., Chen, Z., Xu, D., and Shang, Y. Learning to parallel: Accelerating diffusion large language models via adaptive parallel decoding. *arXiv:2509.25188*, 2025.

Ben-Hamu, H., Gat, I., Severo, D., Nolte, N., and Karrer, B. Accelerated sampling from masked diffusion models via entropy bounded unmasking. *NeurIPS*, 2025.

Bengio, E., Bacon, P.-L., Pineau, J., and Precup, D. Conditional computation in neural networks for faster models. *arXiv:1511.06297*, 2015.

Campbell, A., De Bortoli, V., Shi, J., and Doucet, A. Self-speculative masked diffusions. *arXiv:2510.03929*, 2025.

Chang, H., Zhang, H., Jiang, L., Liu, C., and Freeman, W. T. Maskgit: Masked generative image transformer. In *CVPR*, 2022.

Chen, M., Tworek, J., Jun, H., Yuan, Q., de Oliveira Pinto, H. P., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter, C., Tillet, P., Such, F. P., Cummings, D., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss, A., Guss, W. H., Nichol, A., Paino, A., Tezak, N., Tang, J., Babuschkin, I., Balaji, S., Jain, S., Saunders, W., Hesse, C., Carr, A. N., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford, A., Knight, M., Brundage, M., Murati, M., Mayer, K., Welinder, P., McGrew, B., Amodei, D., McCandlish, S., Sutskever, I., and Zaremba, W. Evaluating large language models trained on code, 2021.

Chen, Z., Fang, G., Ma, X., Yu, R., and Wang, X. dparallel: Learnable parallel decoding for dllms. *arXiv:2509.26488*, 2025.

Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. Training verifiers to solve math word problems. *arXiv:2110.14168*, 2021.

Dai, M., Yang, C., and Si, Q. S-grpo: Early exit via reinforcement learning in reasoning models. *NeurIPS*, 2025.

Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019.

Dieleman, S., Sartran, L., Roshannai, A., Savinov, N., Ganin, Y., Richemond, P. H., Doucet, A., Strudel, R., Dyer, C., Durkan, C., et al. Continuous diffusion for categorical data. *arXiv:2211.15089*, 2022.

Dong, Y., Ma, Z., Jiang, X., Fan, Z., Qian, J., Li, Y., Xiao, J., Jin, Z., Cao, R., Li, B., et al. Saber: An efficient sampling with adaptive acceleration and backtracking enhanced re-masking for diffusion language model. *arXiv:2510.18165*, 2025.

Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al. The llama 3 herd of models. *arXiv e-prints*, 2024.

Gong, S., Zhang, R., Zheng, H., Gu, J., Jaitly, N., Kong, L., and Zhang, Y. Diffucoder: Understanding and improving masked diffusion models for code generation. *arXiv:2506.20639*, 2025.

Graves, A. Adaptive computation time for recurrent neural networks. *arXiv:1603.08983*, 2016.

Guo, G. and Ermon, S. Reviving any-subset autoregressive models with principled parallel sampling and speculative decoding. *arXiv:2504.20456*, 2025.

- 495 Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart,
496 S., Tang, E., Song, D., and Steinhardt, J. Measuring math-
497 ematical problem solving with the math dataset. *NeurIPS*,
498 2021.
- 499 Hong, C., An, S., Kim, M.-S., and Ye, J. C. Improving
500 discrete diffusion unmasking policies beyond explicit
501 reference policies. *arXiv:2510.05725*, 2025a.
- 503 Hong, F., Yu, G., Ye, Y., Huang, H., Zheng, H., Zhang,
504 Y., Wang, Y., and Yao, J. Wide-in, narrow-out:
505 Revokable decoding for efficient and effective dlms.
506 *arXiv:2507.18578*, 2025b.
- 507 Hoogetboom, E., Nielsen, D., Jaini, P., Forré, P., and Welling,
508 M. Argmax flows and multinomial diffusion: Learning
509 categorical distributions. *NeurIPS*, 2021.
- 511 Huang, P., Liu, S., Liu, Z., Yan, Y., Wang, S., Chen, Z., and
512 Xiao, T. Pc-sampler: Position-aware calibration of decod-
513 ing bias in masked diffusion models. *arXiv:2508.13021*,
514 2025a.
- 516 Huang, Z., Chen, Z., Wang, Z., Li, T., and Qi, G.-J. Reinforc-
517 ing the diffusion chain of lateral thought with diffusion
518 language models. *NeurIPS*, 2025b.
- 519 Huang, Z., Wang, Y., Chen, Z., and Qi, G.-J. Don't settle too
520 early: Self-reflective remasking for diffusion language
521 models. *arXiv preprint arXiv:2509.23653*, 2025c.
- 523 Israel, D., Broeck, G. V. d., and Grover, A. Accelerating
524 diffusion llms via adaptive parallel decoding. *NeurIPS*,
525 2025.
- 526 Jiang, Y., Cai, Y., Luo, X., Fu, J., Wang, J., Liu, C., and
527 Yang, X. d^2 cache: Accelerating diffusion-based llms via
528 dual adaptive caching. *arXiv:2509.23094*, 2025.
- 530 Kim, J., Shah, K., Kontonis, V., Kakade, S., and Chen,
531 S. Train for the worst, plan for the best: Understanding
532 token ordering in masked diffusions. *ICML*, 2025a.
- 533 Kim, S. H., Hong, S., Jung, H., Park, Y., and Yun, S.-
534 Y. Klass: Kl-guided fast inference in masked diffusion
535 models. *NeurIPS*, 2025b.
- 537 Kingma, D. P. and Gao, R. Understanding diffusion objec-
538 tives as the ELBO with simple data augmentation. In
539 *NeurIPS*, 2023.
- 541 Li, J., Dong, X., Zang, Y., Cao, Y., Wang, J., and Lin, D.
542 Beyond fixed: Training-free variable-length denoising
543 for diffusion large language models. *arXiv:2508.00819*,
544 2025a.
- 545 Li, P., Zhou, Y., Muhtar, D., Yin, L., Yan, S., Shen, L., Liang,
546 Y., Vosoughi, S., and Liu, S. Diffusion language models
547 know the answer before decoding. *arXiv:2508.19982*,
548 2025b.
- 549 Lin, N., Zhang, J., Hou, L., and Li, J. Boundary-guided
policy optimization for memory-efficient rl of diffusion
large language models. *arXiv:2510.11683*, 2025.
- Liu, S., Nam, J., Campbell, A., Stärk, H., Xu, Y., Jaakkola,
T., and Gómez-Bombarelli, R. Think while you generate:
Discrete diffusion with planned denoising. *ICLR*, 2025.
- Lou, A., Meng, C., and Ermon, S. Discrete diffusion mod-
eling by estimating the ratios of the data distribution. In
ICML, 2024.
- Luce, R. D. *Individual choice behavior*, volume 4. Wiley
New York, 1959.
- Ma, Y., Du, L., Wei, L., Chen, K., Xu, Q., Wang, K.,
Feng, G., Lu, G., Liu, L., Qi, X., et al. dinfer: An effi-
cient inference framework for diffusion language models.
arXiv:2510.08666, 2025.
- Nie, S., Zhu, F., You, Z., Zhang, X., Ou, J., Hu, J., Zhou, J.,
Lin, Y., Wen, J.-R., and Li, C. Large language diffusion
models. *NeurIPS*, 2025.
- Ou, J., Nie, S., Xue, K., Zhu, F., Sun, J., Li, Z., and Li,
C. Your absorbing discrete diffusion secretly models the
conditional distributions of clean data. *ICLR*, 2025.
- Plackett, R. L. The analysis of permutations. *Journal of the
Royal Statistical Society Series C: Applied Statistics*, 24
(2), 1975.
- Radford, A., Narasimhan, K., Salimans, T., Sutskever, I.,
et al. Improving language understanding by generative
pre-training. 2018.
- Rojas, K., Lin, J., Rasul, K., Schneider, A., Nevmyvaka, Y.,
Tao, M., and Deng, W. Improving reasoning for diffusion
language models via group diffusion policy optimization.
arXiv:2510.08554, 2025.
- Sahoo, S., Arriola, M., Schiff, Y., Gokaslan, A., Marroquin,
E., Chiu, J., Rush, A., and Kuleshov, V. Simple and
effective masked diffusion language models. *NeurIPS*,
2024.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and
Klimov, O. Proximal policy optimization algorithms.
arXiv:1707.06347, 2017.
- Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang,
H., Zhang, M., Li, Y., Wu, Y., et al. Deepseekmath:
Pushing the limits of mathematical reasoning in open
language models. *arXiv:2402.03300*, 2024.
- Shen, J., Sarkar, G., Ro, Y., Sridhar, S. N., Wang, Z.,
Akella, A., and Kundu, S. Improving the throughput of
diffusion-based large language models via a training-free
confidence-aware calibration. *arXiv:2512.07173*, 2025.

- 550 Shi, J., Han, K., Wang, Z., Doucet, A., and Titsias, M. K.
551 Simplified and generalized masked diffusion for discrete
552 data. In *NeurIPS*, 2024.
- 553 Song, Y., Zhang, Z., Luo, C., Gao, P., Xia, F., Luo, H.,
554 Li, Z., Yang, Y., Yu, H., Qu, X., et al. Seed diffusion:
555 A large-scale diffusion language model with high-speed
556 inference. *arXiv:2508.02193*, 2025.
- 557 Sutton, R. S., Barto, A. G., et al. *Reinforcement learning:
558 An introduction*. MIT press Cambridge, 1998.
- 559 Swerdlow, A., Prabhudesai, M., Gandhi, S., Pathak, D., and
560 Fragkiadaki, K. Unified multimodal discrete diffusion.
561 *arXiv:2503.20853*, 2025.
- 562 Tang, X., Dolga, R., Yoon, S., and Bogunovic, I. wd1:
563 Weighted policy optimization for reasoning in diffusion
564 language models. *arXiv:2507.08838*, 2025.
- 565 Teerapittayanon, S., McDanel, B., and Kung, H.-T.
566 Branchynet: Fast inference via early exiting from deep
567 neural networks. In *International conference on pattern
568 recognition (ICPR)*. IEEE, 2016.
- 569 Wang, C., Rashidinejad, P., Su, D., Jiang, S., Wang, S., Zhao,
570 S., Zhou, C., Shen, S. Z., Chen, F., Jaakkola, T., et al.
571 Spg: Sandwiched policy gradient for masked diffusion
572 language models. *arXiv:2510.09541*, 2025a.
- 573 Wang, G., Schiff, Y., Sahoo, S. S., and Kuleshov, V. Re-
574 masking discrete diffusion models with inference-time
575 scaling. *NeurIPS*, 2025b.
- 576 Wang, W., Fang, B., Jing, C., Shen, Y., Shen, Y., Wang, Q.,
577 Ouyang, H., Chen, H., and Shen, C. Time is a feature: Ex-
578 ploiting temporal dynamics in diffusion language models.
579 *arXiv:2508.09138*, 2025c.
- 580 Wang, X., Yu, F., Dou, Z.-Y., Darrell, T., and Gonzalez,
581 J. E. Skipnet: Learning dynamic routing in convolutional
582 networks. In *ECCV*, 2018.
- 583 Wang, X., Xu, C., Jin, Y., Jin, J., Zhang, H., and Deng, Z.
584 Diffusion llms can do faster-than-ar inference via discrete
585 diffusion forcing. *arXiv:2508.09192*, 2025d.
- 586 Wang, Y., Yang, L., Li, B., Tian, Y., Shen, K., and Wang,
587 M. Revolutionizing reinforcement learning framework
588 for diffusion large language models. *arXiv:2509.06949*,
589 2025e.
- 590 Wei, Q., Zhang, Y., Liu, Z., Liu, D., and Zhang, L. Acceler-
591 ating diffusion large language models with slowfast: The
592 three golden principles. *arXiv:2506.10848*, 2025.
- 593 Wu, C., Zhang, H., Xue, S., Liu, Z., Diao, S., Zhu, L.,
594 Luo, P., Han, S., and Xie, E. Fast-dllm: Training-free
595 acceleration of diffusion llm by enabling kv cache and
596 parallel decoding. *arXiv:2505.22618*, 2025.
- 597 Wu, S. and Zhang, J. Free draft-and-verification: Toward
598 lossless parallel decoding for diffusion large language
599 models. *arXiv:2510.00294*, 2025.
- 600 Xu, Z., Liu, Y., Yin, Y., Zhou, M., and Poovendran, R.
601 Kodcode: A diverse, challenging, and verifiable synthetic
602 dataset for coding. *arXiv:2503.02951*, 2025.
- 603 Ye, J., Xie, Z., Zheng, L., Gao, J., Wu, Z., Jiang, X., Li,
604 Z., and Kong, L. Dream 7b: Diffusion large language
605 models. *arXiv:2508.15487*, 2025.
- 606 Yi, J., Wang, J., and Li, S. Shorterbetter: Guiding reason-
607 ing models to find optimal inference length for efficient
608 reasoning. *NeurIPS*, 2025.
- 609 Yu, R., Ma, X., and Wang, X. Dimple: Discrete diffusion
610 multimodal large language model with parallel decoding.
611 *arXiv:2505.16990*, 2025.
- 612 Zhan, A. Principled and tractable rl for reasoning with
613 diffusion language models. *arXiv:2510.04019*, 2025.
- 614 Zhao, H., Liang, D., Tang, W., Yao, D., and Kallus, N.
615 Diffpo: Training diffusion llms to reason fast and furious
616 via reinforcement learning. *arXiv:2510.02212*, 2025a.
- 617 Zhao, S., Gupta, D., Zheng, Q., and Grover, A. d1: Scaling
618 reasoning in diffusion large language models via rein-
619 forcement learning. *NeurIPS*, 2025b.
- 620 Zhou, R., Ni, Z., Chen, T., Liu, Z., Yue, Y., Wang, Y., Wang,
621 Y., Liu, J., and Huang, G. Co-grpo: Co-optimized group
622 relative policy optimization for masked diffusion model.
623 *arXiv:2512.22288*, 2025.
- 624 Zhu, Y., Guo, W., Choi, J., Molodyk, P., Yuan, B., Tao,
625 M., and Chen, Y. Enhancing reasoning for diffu-
626 sion llms via distribution matching policy optimization.
627 *arXiv:2510.08233*, 2025.

Appendix

The appendix is organized as follows:

- In [Appendix A](#), we describe related work.
- In [Appendix B](#), we provide a diagram and an algorithm for our proposed policy sampling.
- In [Appendix C](#), we present additional plots to supplement the experiments from the main paper:
 - In [C.1](#), we replicate [Figure 1](#) across models (LLaDA, Dream) and datasets (GSM8k, MATH).
 - In [C.2](#), we replicate [Figure 3](#) using latency (wall-clock time) as the efficiency measure.
 - In [C.3](#), we replicate [Figure 3](#) for policies trained on Dream.
 - In [C.4](#), we show results for the same setting as [Figure 3a](#) and [Figure 3c](#) but with a denser α -grid. We also show the spread of policies across different training seeds.
 - In [C.5](#), we plot the impact of varying the policy temperature τ_π .
 - In [C.6](#), we provide model transfer results for the MATH dataset.
 - In [C.7](#), we provide generation length (L) transfer results for the MATH dataset.
- In [Appendix D](#), we perform a qualitative analysis of our learned sampling policies.
- In [Appendix E](#), we visualize token generation trajectories.
- In [Appendix F](#), we describe dynamic Plackett-Luce sampling as an alternative to Bernoulli sampling.
- In [Appendix G](#), we detail our *expert steering* approach.
- In [Appendix H](#), we discuss an extended background on discrete diffusion models.
- In [Appendix I](#), we provide implementation and hyperparameters details.
- In [Appendix J](#), we visualize the architecture of our sampling policy.

A. Related work

Heuristic Samplers for dLLMs. Sampling in diffusion LLMs (dLLMs) (Nie et al., 2025; Ye et al., 2025) has recently attracted significant attention, with much of the work in this area proposing heuristic approaches to improve the decoding process in a training-free manner. Throughout this paper, we focus on Fast-dLLM (Wu et al., 2025) as a representative heuristic method, as it popularized confidence-thresholded sampling in dLLMs and demonstrated its crucial role in enabling faster inference in dLLMs compared to autoregressive models. Many recently proposed heuristics can be viewed as either reinterpretations or extensions of such confidence thresholding strategies (Ben-Hamu et al., 2025; Wei et al., 2025; Hong et al., 2025b; Li et al., 2025b; Yu et al., 2025; Kim et al., 2025b; Shen et al., 2025). Other notable heuristic approaches explore incorporating spatial (Huang et al., 2025a) or temporal information (Wang et al., 2025c), alternative confidence measures (Kim et al., 2025a), or more explicit modeling of token dependencies (Azangulov et al., 2025). In this work, we aim to complement ongoing research on sampling heuristics by investigating whether effective sampling strategies can be learned directly via reinforcement learning. Note that beyond improving the efficiency of unmasking in dLLMs, heuristics have also been proposed for remasking (Hong et al., 2025b; Dong et al., 2025) and for dynamically adjusting the generation length (Li et al., 2025a), which we do not target in this work.

Reinforcement Learning Post-Training for dLLMs. Early work on post-training diffusion LLMs via reinforcement learning includes d1 (Zhao et al., 2025b), which introduces a variant of GRPO tailored to dLLMs, and DiffuCoder (Gong et al., 2025), which focuses on enhancing the coding abilities of LLaDA-style models using RL. Most recent extensions aim to improve the quality of policy gradient estimators (Tang et al., 2025; Wang et al., 2025a; Lin et al., 2025; Rojas et al., 2025; Zhu et al., 2025; Wang et al., 2025e; Zhan, 2025), and have demonstrated promising results in further enhancing the reasoning capabilities of dLLMs. The key distinction from our approach is that these methods use a fixed sampling strategy (e.g., high-confidence sampling), and the policy corresponds to the dLLM itself (or a LoRA-augmented version for efficiency). Closest to our work are DCOLT (Huang et al., 2025b), which trains a separate unmasking module (with a fixed number of unmasked tokens in each step) in addition to updating the base model via RL, and DiFFPO (Zhao et al., 2025a), which jointly learns an unmasking confidence threshold while updating the base model through RL. Differently to both these works, our primary goal is not in improving the reasoning abilities of the underlying dLLM; hence we keep the underlying dLLM fixed and focus solely on training a standalone policy with the aim of learning fast and adaptive sampling while preserving the performance of the base model. Another related concurrent work is Hong et al. (2025a) where an unmasking policy for dLLMs trained via GRPO is proposed, similar to our approach. However, their policy unmask a fixed number of tokens at each step (similar to DCOLT), whereas ours dynamically adapts the number of tokens to unmask via our Bernoulli formulation. Finally, Seed Diffusion (Song et al., 2025) identifies computation-aware reinforcement learning as a key ingredient for achieving competitive efficiency in closed-source coding dLLMs.

Orthogonal efforts to accelerate dLLMs. Besides the aforementioned heuristic-based sampling innovations, other efforts to improve inference efficiency in dLLMs (Ma et al., 2025) include KV caching (Jiang et al., 2025; Wu et al., 2025), (variants of) speculative decoding (Israel et al., 2025; Campbell et al., 2025; Guo & Ermon, 2025; Wu & Zhang, 2025), training separate decoder modules during pretraining (Liu et al., 2025; Arriola et al., 2025b), and diffusion forcing (Wang et al., 2025d), among others. Concurrent work explores distilling faster generation patterns directly into the base model (Chen et al., 2025), or training a (small) separate unmasking network on top of the pretrained dLLM (Bansal & Sanghavi, 2025; Bao et al., 2025). Crucially, unlike our approach, where this training is done via RL, these works train unmasking modules by distilling generation trajectories from the base model. As a result, we expect that such approaches may still inherit limitations observed in dLLM sampling, such as difficulties in the non semi-AR regime (cf. Section 2.2).

RL for Adaptive Compute Since our sampling policies result in a variable number of sampling steps ($T - \hat{T}$) per input, our work connects to the broader literature on adaptive computation (Graves, 2016; Teerapittayanon et al., 2016). Notable examples of using RL to learn input-dependent compute policies include conditional computation with stochastic gating (Bengio et al., 2015), dynamic block skipping in residual networks (Wang et al., 2018), and RL-based early-exit decisions for long chain-of-thought reasoning (Dai et al., 2025). To the best of our knowledge, our work is the first to explore learning adaptive policies using RL for the task of sampling in dLLMs. The closest related concurrent work is DiFFPO (Zhao et al., 2025a); however, unlike our approach—where the policy is learned end-to-end—DiFFPO learns to predict adaptive thresholds, which are then used in the same fashion as the fixed thresholds employed by Fast-dLLM (Wu et al., 2025). Concurrently, there is also a growing body of work that uses RL to introduce adaptivity into reasoning chains of autoregressive LLMs (Arora & Zanette, 2025; Yi et al., 2025).

B. Policy Sampling Diagram and Algorithm

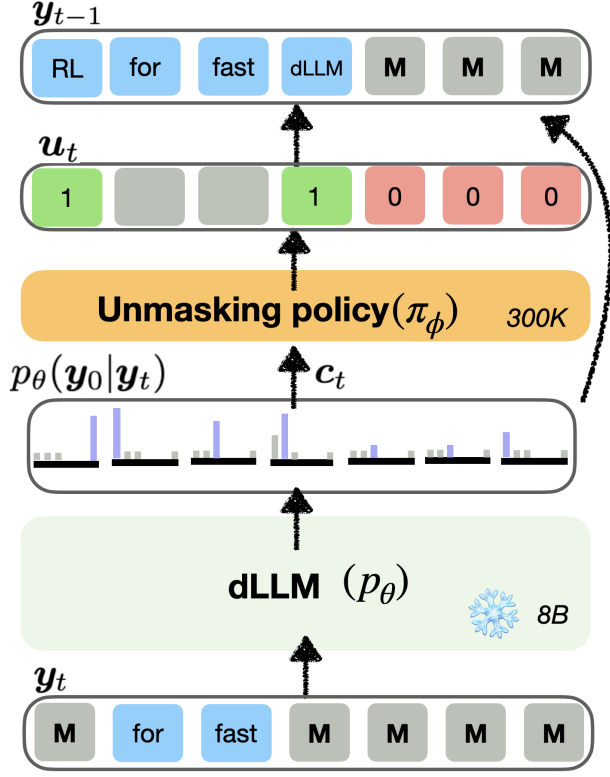


Figure 8. Diagram showing how our policies are used on top of a pretrained dLLM to unmask tokens and generate text.

Algorithm 1. Policy Sampling (full-diffusion setting $BL = L$)

```

1: Input: Prompt  $x \in \mathcal{V}^d$ , maximum generation length  $L$ , maximum diffusion steps  $T$ , dLLM  $p_\theta$ , sampling policy  $\pi_\phi$  (via  $f_\phi$ ), policy temperature  $\tau_\pi$ 
2: Output: Generated sequence  $y_{\hat{T}}$ , number of sampling steps
3:  $y_T \leftarrow (M, \dots, M)$   $\triangleright$  Fully masked initial sequence
4:  $\mathcal{M}_T \leftarrow \{1, \dots, L\}$   $\triangleright$  Set of still-masked positions
5: for  $t = T, T - 1, \dots, 1$  do
6:   if  $\mathcal{M}_t = \emptyset$  then
7:     break  $\triangleright$  Stop generating when all unmasked
8:   end if
9:    $p_t^k \leftarrow p_\theta(\cdot | x, y_t), \forall k \in [L]$   $\triangleright$  Token distributions
10:   $c_t^k \leftarrow \max_{v \in \mathcal{V}} p_t^k(v), \forall k \in [L]$   $\triangleright$  Token confidences
11:   $c_t \leftarrow (c_t^1, \dots, c_t^L)$ 
12:   $m_t^k \leftarrow \mathbf{1}[k \in \mathcal{M}_t], \forall k \in [L]$ 
13:   $b_t \leftarrow f_\phi(c_t, m_t, t)$   $\triangleright$  Policy logits
14:   $s_t^k \leftarrow \sigma(b_t^k / \tau_\pi), \forall k \in \mathcal{M}_t$ 
15:   $u_t^k \sim \text{Ber}(s_t^k), \forall k \in \mathcal{M}_t$   $\triangleright$  Unmasking decisions
16:   $u_t^k \leftarrow 0, \forall k \notin \mathcal{M}_t$   $\triangleright$  Ignore already-unmasked
17:  if  $\sum_{k=1}^L u_t^k = 0$  then
18:     $u_t^{k^*} \leftarrow 1, k^* = \text{argmax}_{k \in \mathcal{M}_t} s_t^k$   $\triangleright$  Argmax fallback
19:  end if
20:   $\mathcal{U}_t^\pi \leftarrow \{k \in \mathcal{M}_t \mid u_t^k = 1\}$ 
21:   $y_{t-1}^k \sim p_t^k, \forall k \in \mathcal{U}_t^\pi$   $\triangleright$  Sample tokens
22:   $y_{t-1}^k \leftarrow y_t^k, \forall k \notin \mathcal{U}_t^\pi$ 
23:   $\mathcal{M}_{t-1} \leftarrow \mathcal{M}_t \setminus \mathcal{U}_t^\pi$ 
24: end for
25:  $\hat{T} \leftarrow t$ 
26: return  $y_{\hat{T}}, T - \hat{T}$ 
    
```

C. Additional Results

C.1. Figure 1 replicated for {LLaDA-8B-INSTRUCT, DREAM-7B-INSTRUCT} × {GSM8k, MATH-500}

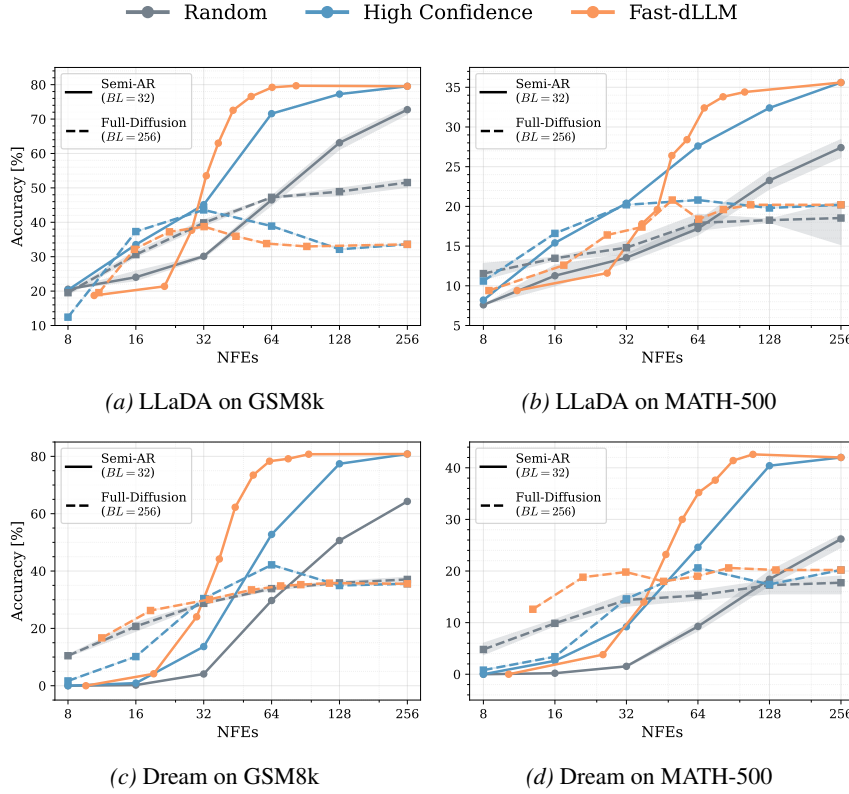


Figure 9. Performance comparison with ($BL = 32$; —) and without ($BL = 256$; - -) semi-AR generation. The same trend observed in Figure 1 holds across all models and datasets: confidence-based heuristics perform well under semi-AR generation but degrade significantly without it.

C.2. Figure 3 replicated using wall-clock time as the efficiency measure

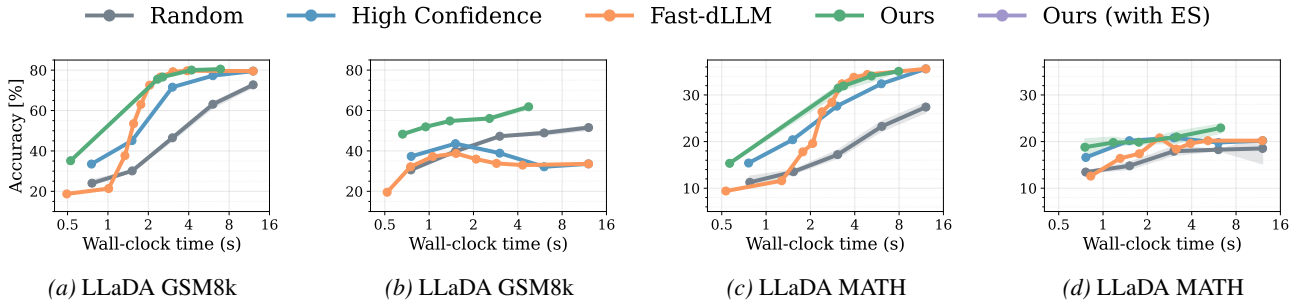


Figure 10. Figure 3 reproduced using wall-clock time (in seconds) as the efficiency metric instead of NFEs. The difference in our policy curves (—) when measured in wall-clock time versus NFEs (see Figure 3) is minimal to non-existent, demonstrating the negligible computational overhead of our policy. This low overhead is primarily due to the small size of the unmasking model relative to the base dLLM (300K vs. 8B parameters in the case of LLaDA). All experiments run on A100 GPUs.

C.3. Figure 3 replicated for DREAM-7B-INSTRUCT

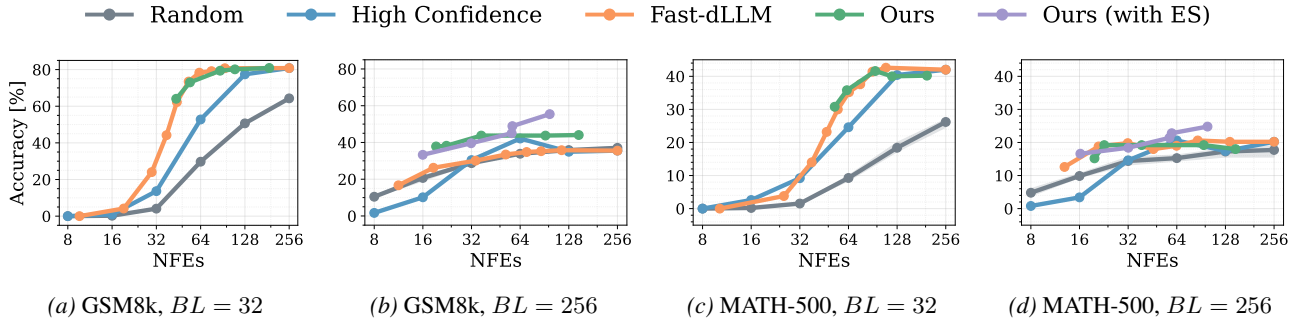


Figure 11. Results for Dream in semi-AR (Figure 11a & Figure 11c) and full-diffusion (Figure 11b & Figure 11d) generation regimes. For the policies we vary $\alpha \in \{10, 3, 1, 0.3, 0\}$ and use $\tau_\pi = 0.5$ for $BL = 32$ and $\tau_\pi = 1$ for $BL = 256$.

C.4. Controllability of learning unmasking policies with RL (via α)

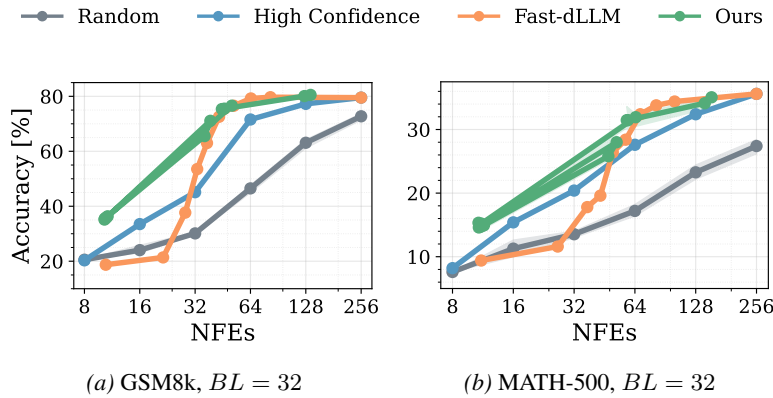


Figure 12. $BL = 32$ results for LLaDA with a denser regularization grid, $\alpha \in \{10.0, 9.0, \dots, 1.0, 0.3, 0.0\}$. Single training seed due to cost; error bars show (min, max) over three test-time seeds. Note that for $\alpha \geq 4.0$, the change in NFEs is not monotonic; different values lead to convergence either close to the $\alpha = 3.0$ policy, or to that of $\alpha = 10.0$. Points for policy results are not connected by α to emphasize the non-monotone behavior.

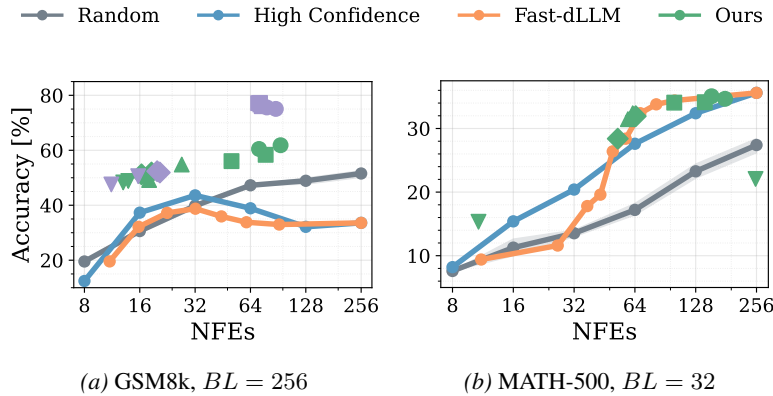


Figure 13. $BL = 256$ (GSM8k) and $BL = 32$ (MATH-500) results for LLaDA with all training seeds scattered. Even for a fixed value of α , the resulting policy can vary in accuracy and speed due to the randomness of the training procedure. Marker shape denotes the value of α .

C.5. Impact of policy temperature τ_π

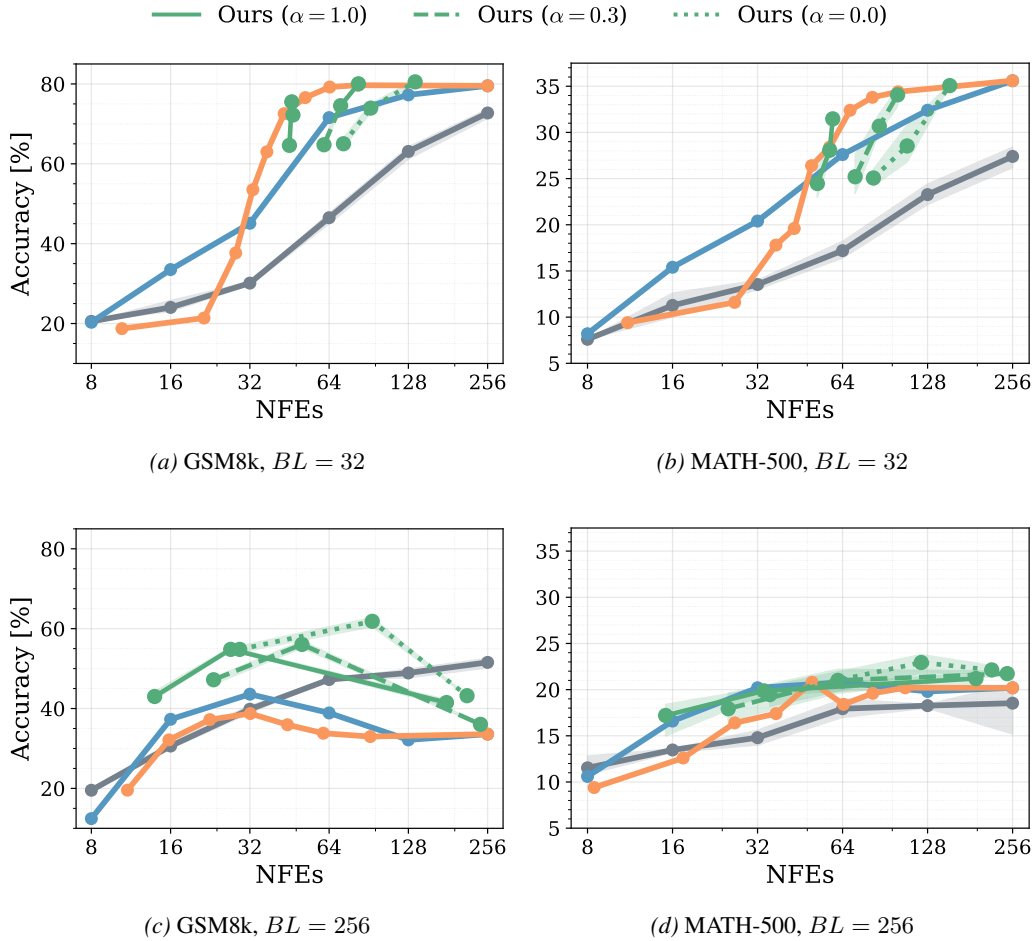


Figure 14. We study the effect of changing the policy temperature τ_π (cf. Section 3.2). For each $\alpha \in \{3, 0.3, 0\}$, we construct a corresponding test-time Pareto frontier by varying $\tau_\pi \in \{1.5, 1.0, 0.5\}$. Interestingly, in some cases—such as $\alpha = 0$ with $BL = 32$ —adjusting τ_π enables an effective trade-off between compute and performance. Moreover, we find that $\tau_\pi = 0.5$ is optimal in the semi-AR, while $\tau_\pi = 1$ performs best in the full-diffusion ($BL = L = 256$) setting.

C.6. Model transfer results

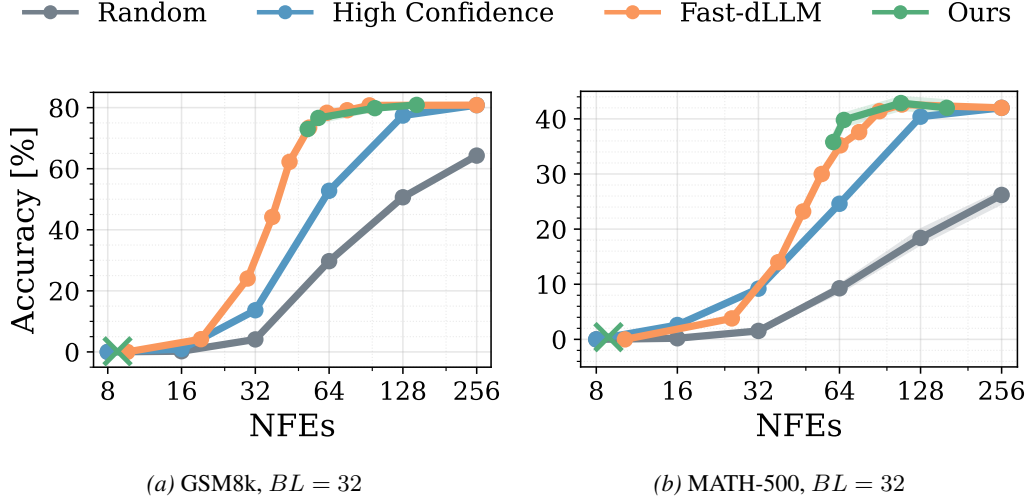


Figure 15. Model transfer results. We use policies trained on LLaDA and evaluate them on Dream with $\tau_\pi = 0.5$. Encouragingly, transferred policies give comparable results to training Dream specific policies (cf. Figure 11a & Figure 11c). Note that the $\alpha = 10$ policy is represented separately by a (x) in the lower left to avoid misleading visualization when interpolating to $\alpha = 3$.

C.7. Sequence-length transferability

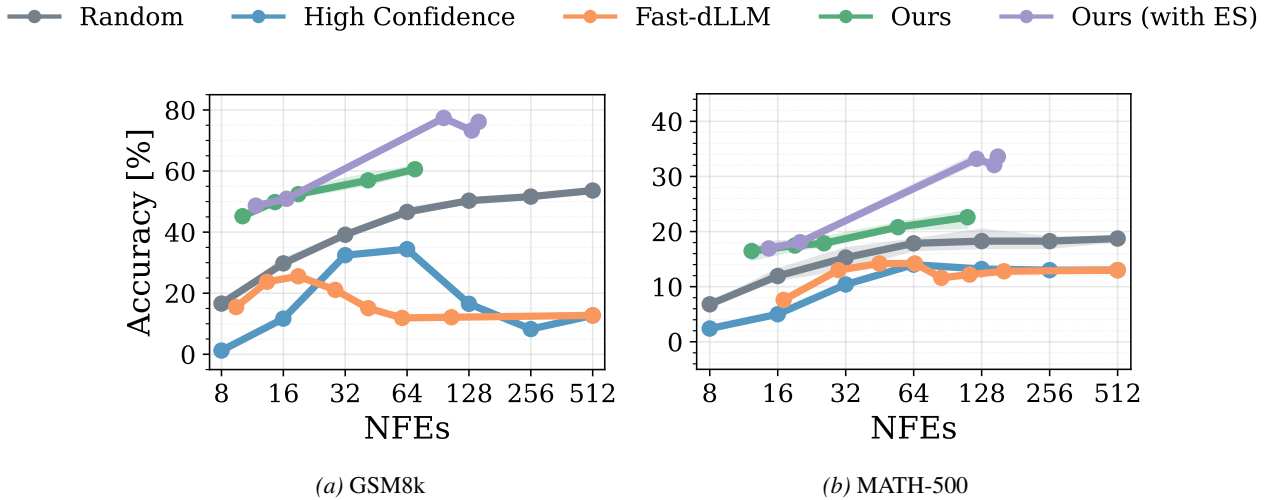


Figure 16. $BL = L = 256$ -trained policies from Section 4.2 evaluated with a 2x longer sequence length ($BL = L = 512$) with $\tau_\pi = 1$. Note that the learned policies yield almost identical performance, while the heuristic methods degrade further compared to $L = 256$ (cf. Figure 3b & Figure 3d). Both for LLaDA-8B-Instruct.

D. What are RL unmasking policies actually doing?

To better understand the behaviour of our learned policies, we conduct a qualitative analysis comparing generation/unmasking trajectories produced by our RL policies against a confidence heuristic like Fast-dLLM (Wu et al., 2025). We consider three settings that span different points on the accuracy-efficiency Pareto frontier and different generation regimes (see Figure 17):

- **BL32-slow** (○): Semi-AR ($BL = 32$) with Fast-dLLM ($\lambda = 0.9$) vs. Policy ($\alpha = 0.3$)
- **BL32-fast** (○): Semi-AR ($BL = 32$) with Fast-dLLM ($\lambda = 0.1$) vs. Policy ($\alpha = 10$)
- **BL256** (○): Full-diffusion ($BL = L = 256$) with Fast-dLLM ($\lambda = 0.9$) vs. Policy ($\alpha = 0.3$) vs. Policy with expert steering ($\alpha = 0.3$)

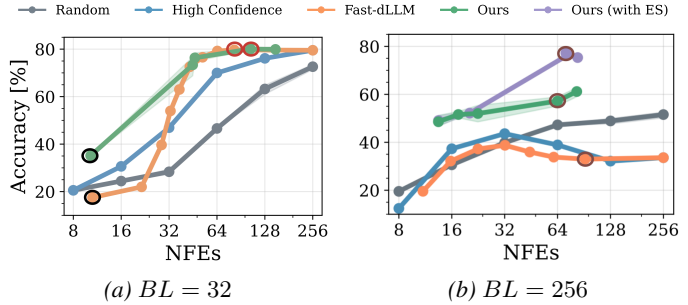


Figure 17. Results for LLaDA in semi-AR ($BL = 32$; figure replicated from Figure 3a) and full-diffusion ($BL = 256$; figure replicated from Figure 3b) generation regimes. We highlight the policies used in our qualitative analysis of unmasking orders: ○ for BL32-slow, ○ for BL32-fast, and ○ for BL256.

Given a prompt x and generated answer \hat{y} (with a maximum answer length and maximum number of steps to set to L and T , respectively), the *unmasking order*² is defined as $\mathbf{o} \in [T]^L$, where $o_k := \min\{t \mid y_t^k \neq M\}$ denotes the timestep when the token at position k is first unmasked. For presentation clarity, we reverse the time axis relative to the main paper, so time flows as $t = 1, 2, \dots, T$. For all experiments, we use $N = 100$ GSM8k test samples with $L = T = 256$.

We describe main insights for each setting next.

D.1. BL32-slow

In this setting, our policy achieves performance nearly identical to Fast-dLLM in both accuracy and efficiency (see ○ in Figure 17a). Despite these similarities, we find that the learned policy adopts a somewhat different sampling strategy compared to plain confidence thresholding. This is visually apparent in Figure 19, where we observe that, although both strategies predominantly follow a left-to-right generation pattern (due to semi-AR regime), the policy’s unmasking order lines (—) display more “wiggleness” within each block. This is due to the fact that, as shown in Figure 20, the policy is far less likely to unmask adjacent tokens simultaneously—doing so for only around $\sim 20\%$ of tokens, in contrast to Fast-dLLM, which does so for $\sim 65\%$ on average. This suggests that the policy learns a more dispersed unmasking strategy, favoring token commitments at scattered positions within each block rather than in contiguous chunks. Such a scattered approach may enable the model to gather more diverse contextual information before finalizing neighboring tokens, potentially mitigating errors that arise from parallel generation. The more scattered generation pattern within each block is also evident in the token generation trajectories shown in Figure 23.

The two methods also differ substantially in how they allocate compute across blocks. As shown in Figure 4a, Fast-dLLM assigns a larger number of NFEs to earlier blocks, allocating ~ 17 NFEs to the first block and gradually decreasing to just ~ 4 NFEs in the final block. In contrast, our learned policy distributes compute almost uniformly across blocks, using ~ 10 NFEs per block. This likely also contributes to the decreasing Spearman rank correlation of unmasking orders between the policy and Fast-dLLM sampling across blocks (see Figure 18).

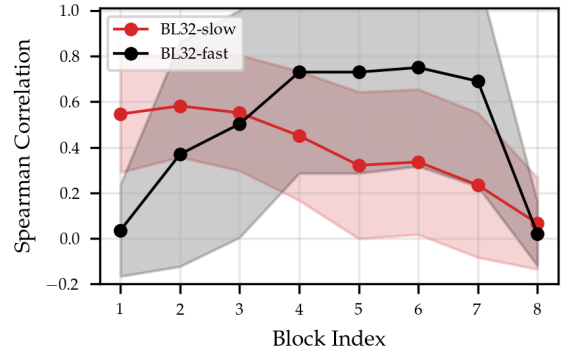


Figure 18. Spearman rank correlation between the unmasking orders of Fast-dLLM and our RL policies per block for semi-AR policies considered here. Note that for certain blocks, the unmasking orders exhibit zero correlation: first and last block for BL32-fast (○), and last block for BL32-slow (○).

²Note that we consider partial orders here, as multiple tokens can be unmasked simultaneously.

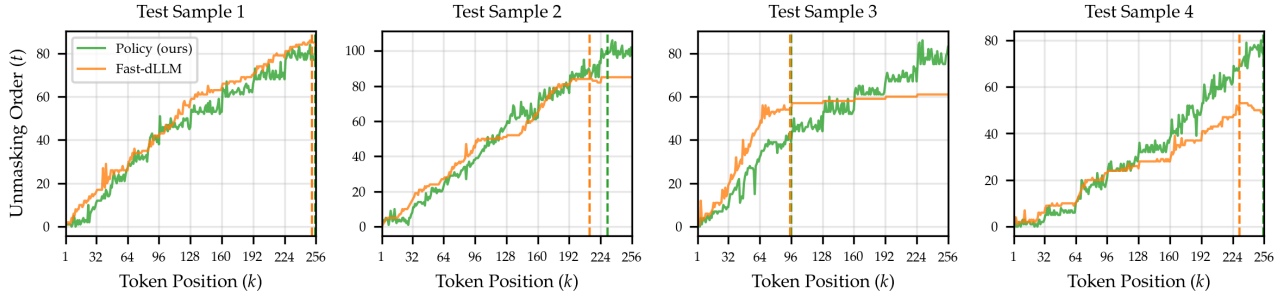


Figure 19. **BL32-slow**: unmasking orders for policy sampling ($\alpha=0.3$) and Fast-dLLM ($\lambda=0.9$) for four GSM8k test samples under semi-AR generation ($BL = 32, L = 256$). Vertical dashed lines indicate the position where the model outputs the end-of-sequence (EOS) token. Interestingly, Fast-dLLM appears more effective at accelerating generation after producing the EOS token, whereas policy sampling tends to sometimes waste compute by not unmasking all padding tokens (within the current block) after EOS in parallel, see the third test sample.

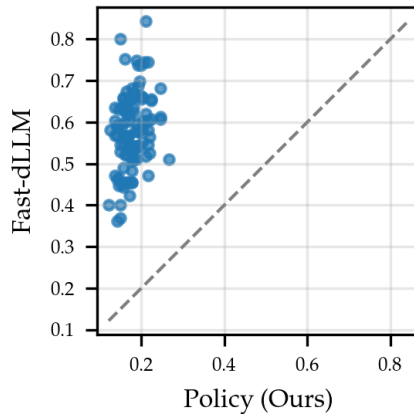


Figure 20. **BL32-slow**: for each of the $N = 100$ samples, we report the frequency of tokens whose right-adjacent token was unmasked at the same sampling step under Fast-dLLM versus our policy sampling. Interestingly, policy sampling un.masks adjacent tokens much less frequently than Fast-dLLM.

D.2. BL32-fast

Here both strategies achieve very fast but low-accuracy sampling (see \bigcirc in Figure 17a). However, policy sampling significantly outperforms Fast-dLLM, achieving approximately $\sim 38\%$ accuracy compared to $\sim 18\%$. Examining the unmasking orders in Figure 21 and the average NFEs per block in Figure 4b, we observe that policy sampling tends to slow down in the final block, typically generating the numerical answer last (see also the token generation trajectory in Figure 25 and how the tokens of the numerical answer are generated last under policy sampling). This behavior may help explain its improved performance over Fast-dLLM and shows the promise of RL for discovering new sampling techniques.

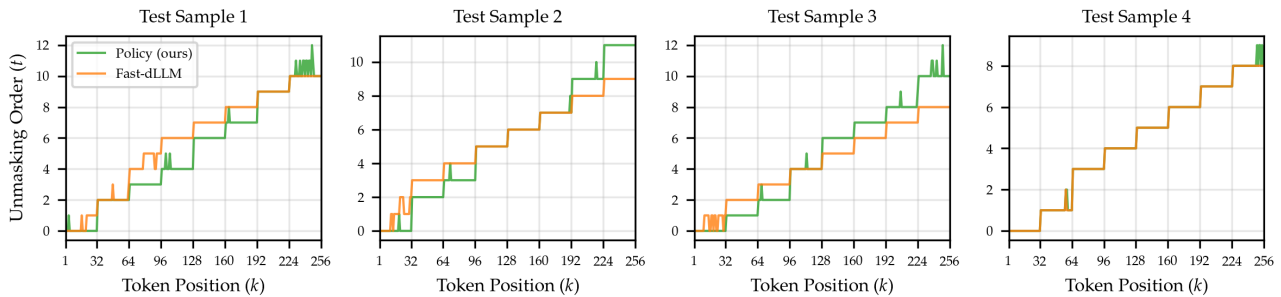


Figure 21. **BL32-fast**: unmasking orders for policy sampling ($\alpha=10$) and Fast-dLLM ($\lambda=0.1$) for four GSM8k test samples under semi-AR generation ($BL = 32, L = 256$). Under low-NFE constraints, both strategies often predict all tokens within a block simultaneously (as indicated by flat unmasking-order lines). In the final block, when generating the numerical answer, the policy slows down.

D.3. BL256

Lastly, we consider the full-diffusion setting ($BL = L = 256$). While all three strategies achieve a comparable number of NFEs (see \circ in Figure 17b), their accuracy differs significantly. The fact that these strategies implement notably different sampling procedures is further evidenced by the low values of Spearman rank correlations 0.26 ± 0.21 when comparing Fast-dLLM orders with policy sampling, and -0.24 ± 0.51 when comparing Fast-dLLM with policy trained with expert steering. These differences are also reflected in the divergent unmasking orders shown in Figure 22 (also refer to Figure 5 for a plot of average unmasking order across samples).

Starting with Fast-dLLM (\rightarrow), we observe in Figure 22 that it often generates tokens in a "reverse" order (right-to-left), see also the token trajectories in Figure 28). This behavior likely stems from the fact that LLaDA did not exclude padding tokens from the loss computation during SFT (Nie et al., 2025) to preserve the dLLM’s ability to generate sequences of varying (effective) lengths. Consequently, the model overfits to padding tokens, which confidence-based methods like Fast-dLLM then mistakenly prioritize, leading to their early generation. This helps explain the poor performance of confidence-based strategies observed in Figure 1, highlighting the fragility of heuristic-driven approaches.

Looking at the unmasking orders of policy sampling (\leftarrow), the policy appears to implement a near-random unmasking strategy, as there is little evident spatial structure in the unmasking order. This is encouraging, as it suggests that the policy can learn to disregard corrupted signals—such as inflated confidence for padding tokens—without explicit supervision, leading to more robust performance across different generation settings (semi-AR vs full-diffusion). Lastly, we observe that the policy with expert steering (\dashrightarrow) generates in an AR fashion (left-to-right). This behavior likely accounts for its superior performance in the full-diffusion setting. It indicates that, during full-diffusion training, the policy is able to pick up the accuracy benefits of AR generation from its expert demonstrations—drawn from the semi-AR regime (Appendix G)—while still learning to generate significantly faster than a purely AR strategy that unmaskes one token at a time.

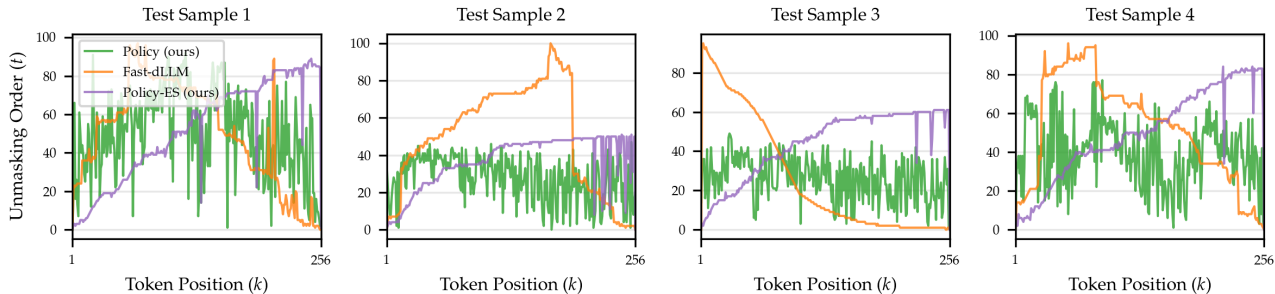


Figure 22. BL256: unmasking orders for policy sampling ($\sigma_\alpha = 0.3$) and Fast-dLLM ($\sigma_{\lambda=0.9}$) for four GSM8k test samples under full-diffusion generation ($BL = L = 256$).

E. Generation Trajectories Visualizations

Here we visualize token generation trajectories for the various sampling strategies considered in our work. Specifically, for a randomly selected GSM8k test sample, we display the generated token at each of the $L = 256$ positions, along with its unmasking time (shown in the bottom right corner of each cell, starting from $t = 0$). For visual clarity, token cells are color-coded according to their unmasking time, with blue indicating earlier and red indicating later unmasking. Note that for semi-AR generations, color-coding is done within each block separately.

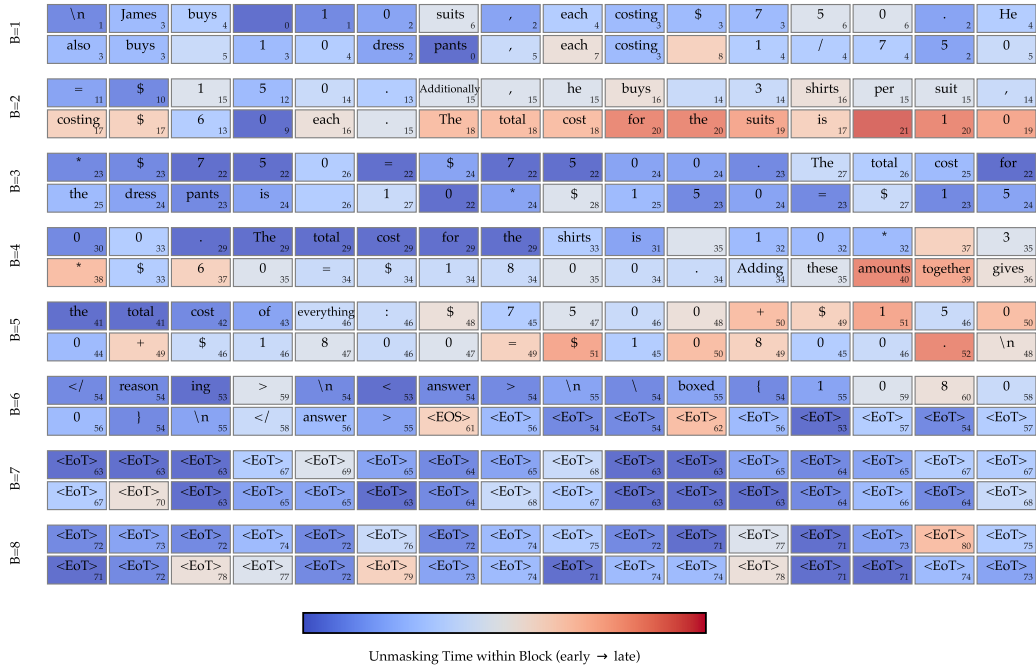


Figure 23. BL32-slow: generation trajectory for policy sampling ($\alpha = 0.3$, semi-AR).

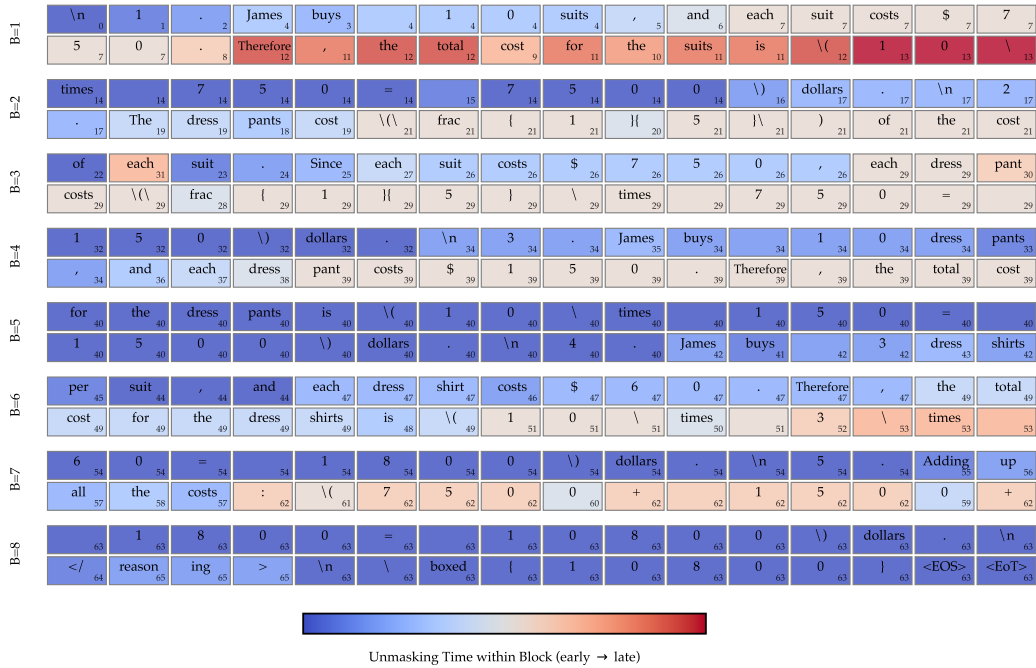


Figure 24. BL32-slow: generation trajectory for Fast-dLLM sampling ($\lambda = 0.9$, semi-AR).

Learning Unmasking Policies for Diffusion Language Models

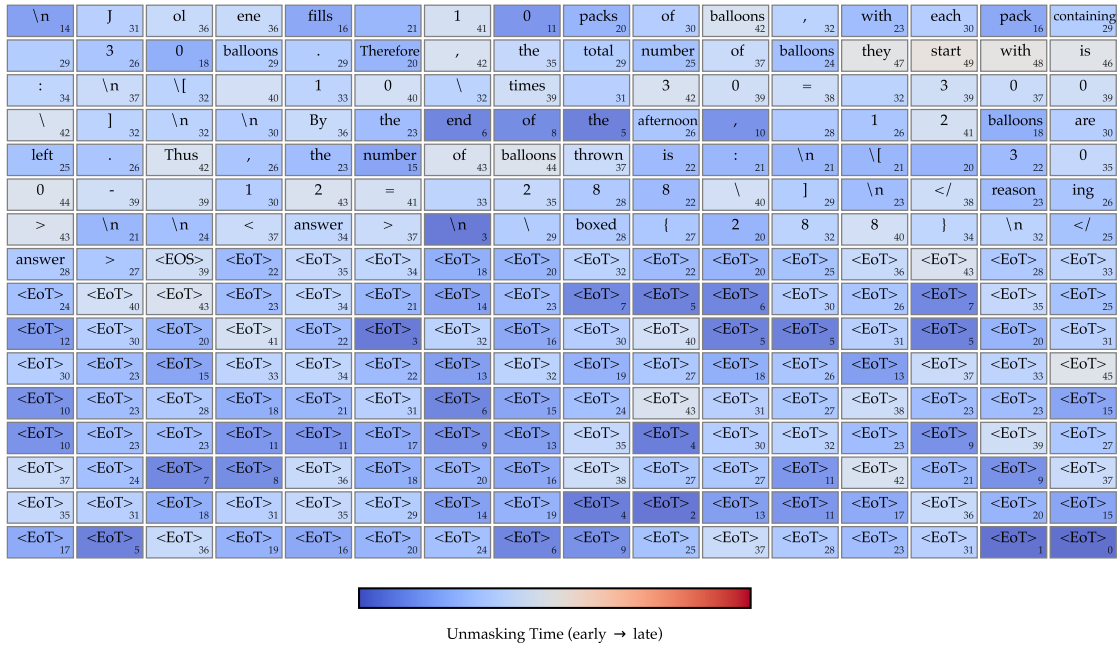


Figure 27. BL256: generation trajectory for policy sampling ($\alpha = 0.3$, full-diffusion).

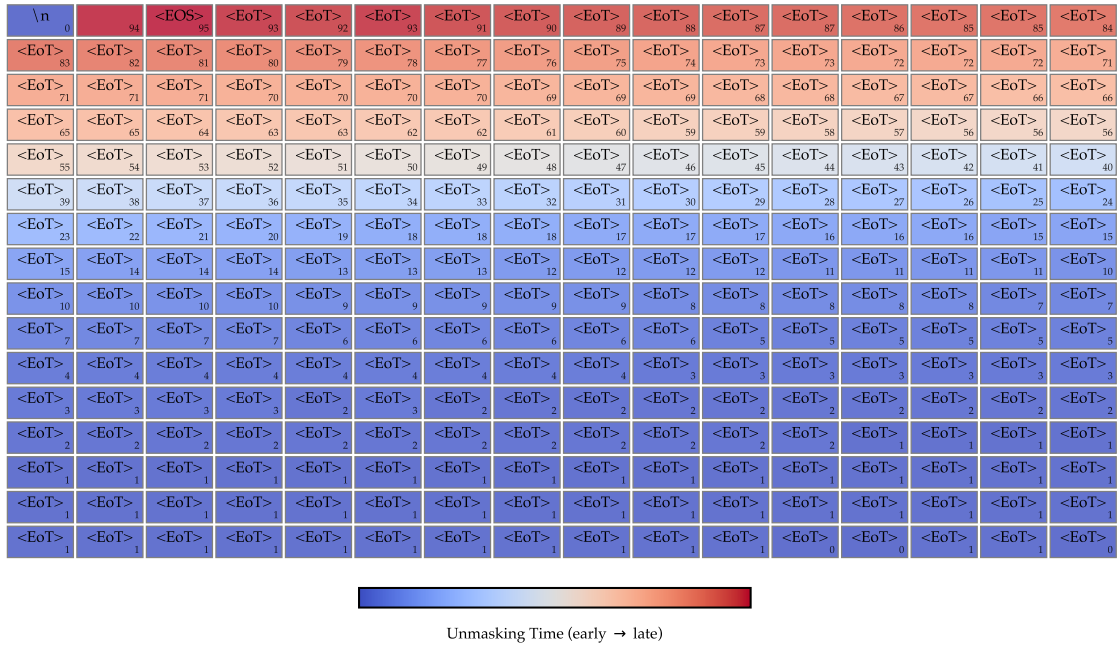
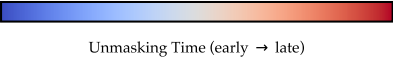


Figure 28. BL256: generation trajectory for Fast-dLLM sampling ($\lambda = 0.9$, full-diffusion).

\n	To	determine	how	many	balloons	were	thrown	,	we	need	to	follow	these	steps	:
\n	1	.	Calculate	the	total	number	of	balloons	initially	.	\n	2	.	Sub	tract
the	number	of	balloons	left	at	the	end	of	the	afternoon	from	the	total	number	of
balloons	.	\n	\n	First	,	we	calculate	the	total	number	of	balloons	:	\n	-
There	are	1	0	packs	of	balloons	.	\n	-	Each	pack	contains	3		
0	balloons	.	\n	-	Therefore	,	the	total	number	of	balloons	is	\n	1	0
\	times	3	0	=	3	0	0	\n).	\n	\n	Next	,		
we	subtract	the	number	of	balloons	left	at	the	end	of	the	afternoon	:	\n	-
There	are	1	2	balloons	left	.	\n	-	Therefore	,	the	number	of	balloons	
thrown	is	\n	3	0	0	-	1	2	=	2	8	8	\n		
).	\n	</	reason	ing	>	\n	<	answer	>	\n	\	boxed	{	2	8
8	}	\n	</	answer	>	<EOS>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>
<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>
<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>
<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>
<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>	<EoT>



Unmasking Time (early → late)

Figure 29. BL256: generation trajectory for policy sampling ($\alpha = 0.3$, expert steering, full-diffusion).

F. Dynamic Plackett-Luce Sampling

We detail here the dynamic Plackett-Luce sampling (DPLS) strategy as an alternative to the Bernoulli sampling (cf. Section 3.2) used in the main experiments of this paper. The name reflects the connection the Plackett-Luce (PL) model proposed by (Luce, 1959) and (Plackett, 1975), with one key difference: the number of selected items is not fixed but can vary freely between 1 and L .

Formally, the Plackett-Luce model operates as follows. Let $\mathbf{b}_t = (b_t^1, \dots, b_t^L)$ denote the unmasking logits (corresponding to the output of policy network f_ϕ , same as in Section 3.2). Interpreting these scores as unnormalized utilities associated with choosing each token for unmasking, under the PL model the likelihood of any particular *permutation* $\sigma := (\sigma_1, \dots, \sigma_L)$ of the tokens is given by

$$P(\sigma \mid \mathbf{b}_t) = \prod_{l \in [L]} \frac{\exp(b_t^{\sigma_l})}{\sum_{j \geq l} \exp(b_t^{\sigma_j})}.$$

Informally, this corresponds to sampling all indices without replacement, where at each step the probability of choosing an item is proportional to its (exponentiated) utility.

The Plackett-Luce model can be easily adapted to model sampling of a fixed-length *ordered* set $\mathcal{U}_t \subseteq [L]$ where $|\mathcal{U}_t| = K$. This is because the probability of any particular partial permutation is equal to the marginalization over all permutations which complete it. Concretely, let $\Sigma(\mathcal{U}_t)$ be the set of permutations which begin with the sequence \mathcal{U}_t , then it follows that

$$P(\mathcal{U}_t \mid \mathbf{b}_t) = \sum_{\sigma \in \Sigma(\mathcal{U}_t)} \underbrace{\prod_{l \in [L]} \frac{\exp(b_t^{\sigma_l})}{\sum_{j \geq l} \exp(b_t^{\sigma_j})}}_{P(\sigma \mid \mathbf{b}_t)} = \prod_{l \in \mathcal{U}_t} \frac{\exp(b_t^{\sigma_l})}{\sum_{j \in \mathcal{U}_t^c \cup \mathcal{U}_t^{\geq l}} \exp(b_t^{\sigma_j})}$$

where \mathcal{U}_t^c denotes the complement of \mathcal{U}_t , and $\mathcal{U}_t^{\geq l}$ denotes the indices which are in \mathcal{U}_t at or after position l .

To extend the PL model to a *variable-length* unmasking sequence, we introduce a special *STOP* token with fixed utility $b_{\text{STOP}} = 0$, and proceed as follows:

1. Sample a single token $l \in [L]$ to unmask

$$l \sim \text{softmax}(\mathbf{b}_t / \tau_\pi)$$

and initialize $\mathcal{U}_t^\pi = \{l\}$.

2. Sample another token l' from the *renormalized distribution*

$$l' \sim \text{softmax}([\mathbf{b}_t^{\setminus \mathcal{U}_t}; b_{\text{STOP}}] / \tau_\pi)$$

where $[\mathbf{b}_t^{\setminus \mathcal{U}_t}; b_{\text{STOP}}]$ denotes \mathbf{b}_t concatenated with b_{STOP} and with the logits of all previously selected indices $l \in \mathcal{U}_t^\pi$ set to $-\infty$.

3. Add l' to \mathcal{U}_t^π .

4. If l' was not *STOP*, repeat from 2.

As written above, this algorithm is not GPU-friendly due to the dynamic computation graph it implies. However, it can be implemented in an efficient manner by treating the loop in steps 2-4 as a Gumbel-argsort and simply masking out all actions which get sampled after *STOP*. Once the samples have been obtained, their likelihoods can be efficiently computed in the same manner as in the fixed-length case above by marginalizing over all actions which occur after *STOP*.

G. Expert Steering

As mentioned in Section 4.2, we find that naively training confidence policies directly on full-diffusion generation ($BL = L$) task yields policies which beat out the heuristic methods, but underperform the ones trained in the semi-AR setting (cf. Figure 3).

Since semi-AR decoding lies within the function class representable by the policy π_ϕ —that is, there is some ϕ such that π_ϕ approximates the semi-AR setting—we hypothesize that our failure to obtain similar policies is due to the vanishingly small probability of encountering autoregressive-like rollouts by chance. More concretely, imagine trying to train a policy on a task for which purely AR generation is substantially more effective than any other decoding strategy. In such a scenario, starting from a randomly initialized policy, the probability of observing a rollout resembling a purely AR sequence is approximately $1/L! \approx 0$ (where L is the generation length), making it extremely unlikely to sample such rollouts during RL training.

To try to eschew this problem, we devise the *Expert Steering* (ES) training strategy as follows. Formally, letting π_ϕ denote the sampling policy to be learned, we simply replace it *at train time only* with the mixture

$$\pi_\phi^{ES}(\cdot | \mathbf{y}_t) = \frac{G}{G + E} \pi_\phi(\cdot | \mathbf{y}_t) + \frac{E}{G + E} \sum_{e=1}^E \delta_e(\cdot | \mathbf{y}_t)$$

where G is the GRPO group size, E is the number of experts to mimic, and each Dirac distribution δ_e represents a chosen deterministic “expert policy” (e.g., a heuristic method). Then in the outer loop of GRPO, instead of sampling a group $\mathbf{g} \sim \pi_\phi$ of size G from π_ϕ , we simply sample an augmented group $\mathbf{g}' \sim \pi_\phi^{ES}$ of size $G + E$. In the inner loop of GRPO we proceed as normal, making sure to use π_ϕ^{ES} when calculating the likelihood ratios to avoid the training instabilities that would otherwise arise (as samples from the Diracs may have likelihood ≈ 0 under π_ϕ)

In practice, for our experiments in Section 4.2 we use a single deterministic expert δ_e corresponding to Fast-dLLM with $\lambda = 0.9$ and $BL = 32$, and force exactly one draw ($E = 1$) from this heuristic per group. The effect is that if the policy is *worse* than this heuristic, that single sample will tend to have positive advantage, causing the policy to be biased toward it. On the other hand, if the policy is *better* than the heuristic, it will have negative advantage, causing it to be made *less* likely. Thus, expert steering encourages exploration towards the expert—in this case, Fast-dLLM—while still allowing the policy to go beyond it, thanks to the group relative loss in GRPO.

H. Extended Background

Masked vs. uniform discrete diffusion. The D3PM framework (Austin et al., 2021a) characterizes forward corruption via a Markov kernel Q_t . Two canonical choices are (i) *masked* (absorbing-state) diffusion, which maps every token toward a distinguished mask state M ,

$$Q_t^{\text{mask}} = (1 - \alpha_t)I + \alpha_t \mathbf{1}e_m^\top,$$

with $\alpha_t \in [0, 1]$, e_m denotes a one-hot vector (with dimension V) for M token, and $\mathbf{1}$ represents a vector of all 1s of size V , and (ii) *uniform-state* diffusion, which replaces tokens uniformly at random,

$$Q_t^{\text{uni}} = (1 - \alpha_t)I + \frac{\alpha_t}{V} \mathbf{1}\mathbf{1}^\top.$$

These differ in their limiting behavior and reverse constraints: Q_t^{mask} has a point-mass stationary distribution at M and induces a monotonic forward trajectory in the mask count (reverse steps deterministically move from $M \rightarrow$ token), whereas Q_t^{uni} has a uniform stationary distribution and permits token \leftrightarrow token substitutions in both directions. The absorbing choice recovers the familiar MDM objective (cf. Equation (1)) under a particular reweighting (Ou et al., 2025) and thus ties masked LMs to discrete diffusion (Austin et al., 2021a).

Continuous-time discrete diffusion. Recent work derives continuous-time limits (CTMC) for discrete diffusion with absorbing corruption and shows that the training objective is a *weighted time integral of token-wise cross-entropy*, with weight proportional to a signal-to-noise ratio term (e.g., $\alpha_t/(1 - \alpha_t)$). This yields a simple, schedule-invariant objective, clarifies forward–reverse consistency, and improves optimization and sampling without changing the prediction target (Shi et al., 2024; Sahoo et al., 2024). Conceptually, this places masked diffusion on the same ODE/SDE footing as continuous-state diffusion (Kingma & Gao, 2023) while retaining native discrete token inference.

Continuous-input diffusion for categorical data. A complementary approach to discrete-state diffusion is to keep the diffusion process continuous in both time and input by operating on token embeddings. For example in (Dieleman et al., 2022), an orthogonal line embeds tokens into \mathbb{R}^d and applies fully continuous diffusion in *both* time and input space. Training can be done with cross-entropy via score interpolation, and sampling uses ODE/SDE solvers with tools like classifier-free guidance. This preserves the continuous machinery but introduces an embedding decoding interface and relaxes exact discreteness during the trajectory.

Summary of differences.

- **Masked (absorbing):** monotonic reverse dynamics ($M \rightarrow$ token only); objective reduces to weighted MDM; pairs naturally with unmask-only sampling policies used in this work.
- **Uniform-state:** non-monotone reverse dynamics (token \leftrightarrow token); encourages broader exploration but typically requires remasking/substitution moves and would thus complicate policy design.
- **Continuous-input:** inherits continuous samplers and guidance; trades exact discreteness for an embedding path and decoding step.

We adopt the absorbing formulation in this work.

I. Training and Policy Network Configuration

Table 1. Training and policy configuration for our main experiments.

Category	Parameter	Value
Training	Learning rate	3e-5
	LR scheduler	Cosine
	Warmup steps	100
	Effective batch size	16
	Weight decay	0.1
	Max gradient norm	0.2
	Clipping factor (ϵ)	0.5 (0.2 for expert steering)
GRPO	Training data	GSM8k and MATH training set mixture
	Num train epochs	1
	Group size	8
	β (KL penalty)	0.0
	Correctness reward $r(\mathbf{y}, \mathbf{y}_t)$	We use d1 (Zhao et al., 2025b) rewards.
Policy Network	Number of transformer blocks	1
	Hidden dimension	128
	Feedforward dimension	512
	Number of heads	2
	Time embedding dim	128
	Total parameter count	300K

J. Policy Architecture Diagram

c = confidence m = mask t = timestep k = top- k d = hidden_dim

