# CreditDecoding: Accelerating Parallel Decoding in Diffusion Large Language Models with Trace Credits

**Kangyu Wang**[1,2*], **Zhiyun Jiang**[1,3,4*], **Haibo Feng**[1], **Weijia Zhao**[1], **Lin Liu**[1],
**Jianguo Li**[1†], **Zhenzhong Lan**[1,4†], **Weiyao Lin**[2†]

[1] Ant Group   [2] Shanghai Jiao Tong University   [3] Zhejiang University   [4] Westlake University
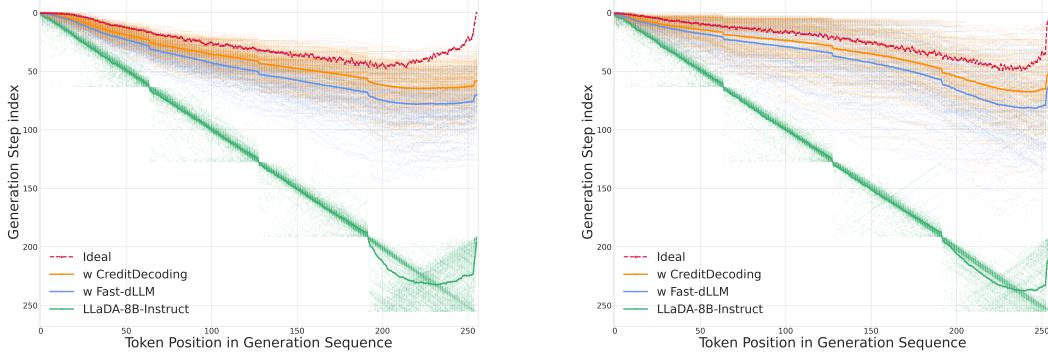
## Abstract

Diffusion large language models (dLLMs) generate text through iterative denoising steps, achieving parallel decoding by denoising only high-confidence positions at each step. However, existing approaches often repetitively remask tokens due to initially low confidence scores, leading to redundant iterations and limiting overall acceleration. Through the analysis of dLLM decoding traces, we observe that the model often determines the final prediction for a token several steps before the decoding step. To leverage this historical information and avoid redundant steps, we introduce the concept of **Trace Credit**, which quantifies each token's convergence potential by accumulating historical logits. Furthermore, we propose **CreditDecoding**, a *training-free* parallel decoding algorithm that accelerates the confidence convergence of correct but underconfident tokens by fusing current logits with Trace Credit. This process significantly reduces redundant iterations and enhances decoding robustness. On eight benchmarks, CreditDecoding achieves a 5.48× speedup and a 0.48 performance improvement over LLaDA-8B-Ins, and a 4.11× speedup with a 0.15 performance improvement over LLaDA-MoE-Ins. Importantly, CreditDecoding scales effectively to long sequences and is orthogonal to mainstream inference optimizations, making it a readily integrable and versatile solution.

## 1 Introduction

Diffusion-based large language models (dLLMs) have recently emerged as a promising alternative to autoregressive (AR) models for text generation (Ye et al., 2025; Nie et al., 2025; Zhu et al., 2025a,b; Gong et al., 2025b,a; Song et al., 2025; Yang et al., 2025; Kim et al., 2025a). Unlike AR models that decode tokens strictly left-to-right, dLLMs generate text through iterative denoising with bidirectional attention, enabling richer contextual dependencies. This paradigm has demonstrated advantages in reasoning and generation quality (Nie et al., 2025; Ye et al., 2025), however, inference efficiency remains a major bottleneck: dLLMs typically require multiple denoising steps to decode all initially masked tokens, and the use of bidirectional attention precludes a lossless KV cache (Yu et al., 2025).

To accelerate inference, recent work has focused on improving the effectiveness and efficiency of *parallel decoding* in dLLMs (Yu et al., 2025; Wei et al., 2025). At each denoising step, the model first predicts all masked tokens and then selects a subset of high-confidence positions to commit, while re-masking the remaining uncertain tokens for future refinement(Yu et al., 2025). This strategy allows multiple tokens to be updated in parallel and has proven both simple and effective. Nevertheless, it suffers from two key limitations: **(i) Computational redundancy.** In many cases, eventually-decoded tokens stabilize early, yet their confidence stay below the threshold and they are re-masked and re-predicted multiple times, wasting compute and limiting parallelism. As illustrated in Figure 1, we quantify this effect via a sizable gap between the first step a token becomes top-1 and

---

*Equal Contribution.   †Corresponding Authors.

| (a) Decoding Boundaries on GSM8K | (b) Decoding Boundaries on HumanEval |

Figure 1: Temporal gap between token stabilization and final decoding. We visualize the *average decoding step* for correct tokens under different parallel decoding strategies on **GSM8K** and **HumanEval**, using **LLaDA-8B-Instruct** with *generation length* $= 256$ and *block size* $= 64$. The *red* curve denotes a theoretical upper bound—the earliest step at which a token is *sampled as the step-wise candidate*. Each other curve shows the *actual* step when that token is finally decoded by the corresponding strategy. The persistent gap between the ideal bound and the actual decoding steps indicates that many correct tokens are repeatedly re-masked after stabilizing, revealing inherent inefficiencies in standard parallel decoding.

the step it is actually decoded. **(ii) History-agnostic Decoding.** During the generation process, as the context continually evolves and may include mispredicted tokens, the confidence of otherwise stable tokens can fluctuate or even regress(Wang et al., 2025). However, token decoding is typically made independently from the current prediction distribution at each step, without leveraging the historical consistency of tokens, which undermines convergence and causes errors to propagate across steps, reducing the robustness of decoding.

To address these issues, we propose **CreditDecoding**, a *training-free* and effective parallel decoding strategy for dLLMs. The key idea is to assign each token a *trace credit score* that accumulates historical logits across steps, and act as a token-level prior that is fused with the current logits to accelerate confidence convergence for correct-but-underconfident tokens. In doing so, CreditDecoding reduces redundant iterations while stabilizing predictions against temporary inconsistencies.

We evaluate CreditDecoding on two dLLMs across eight benchmarks covering knowledge, reasoning, and coding tasks. Our experiments show that CreditDecoding achieves consistent speedups (up to $5.48\times$ in Tokens-Per-Forward) with minimal or no accuracy loss. Furthermore, CreditDecoding is fully orthogonal to existing inference optimizations such as KV cache (Feng et al., 2025) , Early stop and PyTorch 2.0 compiler infrastructure (Jain et al., 2023), making it readily integrable into existing dLLM pipelines.

In summary, our work makes the following contributions:

- We empirically reveal and analyze two key limitations of current dLLM inference under parallel decoding: *computational redundancy* and *history-agnostic decisions*. Through case studies, we reveal a temporal gap between sampling and decoding in dLLM generation, which points to a potential direction for further acceleration.
- We propose **CreditDecoding**, a *training-free* parallel decoding algorithm that accumulates trace credit as a token-level prior to *enhance* and *correct* current logits, thereby accelerating convergence. We demonstrate that on eight benchmarks, CreditDecoding achieves a 5.48× speedup and a 0.48 performance improvement over LLaDA-8B-Ins, and a 4.11× speedup with a 0.15 performance improvement over LLaDA-MoE-Ins..
- Our CreditDecoding is plug-and-play and exhibits scaling potential, it accelerates diverse dLLM backbones and remains effective in long-context (up to 4K) settings. Moreover, it is compatible with mainstream inference optimizations without retraining or architectural changes.
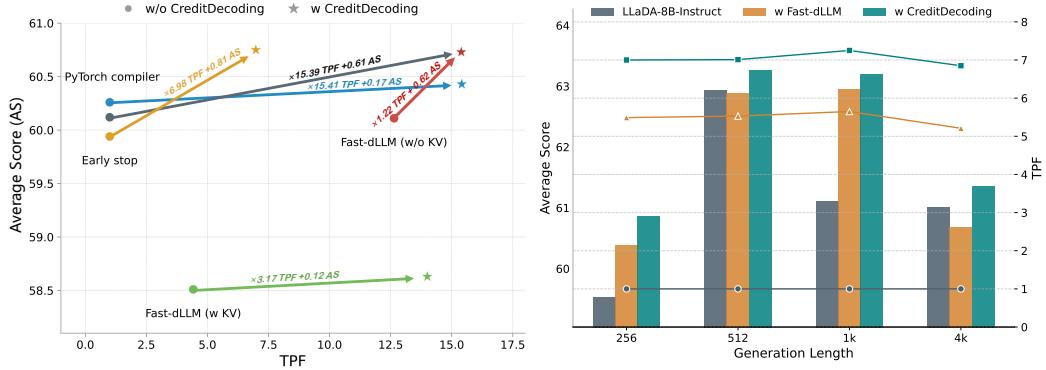
Figure 2: Orthogonality and scalability of CreditDecoding on LLaDA-8B-Instruct. In the left figure(**w/o Early Stop**), each colored dot denotes an acceleration method, and the corresponding star shows its performance with CreditDecoding. In the right figure(**w/ Early Stop**), the line indicates the TPF, while the bars represent the average score.

## 2 Related Work

**Diffusion Language Models**  Diffusion large language models (dLLMs) replace strict left-to-right decoding with iterative denoising, enabling order-agnostic and parallel token updates with bidirectional context (Nie et al., 2025; Ye et al., 2025). Representative systems span general LMs and multimodal variants, including Dream and the LLaDA family (Ye et al., 2025; Nie et al., 2025; Zhu et al., 2025a,b). Recent efforts further scale or specialize dLLMs for coding and large-scale training (Gong et al., 2025b,a; Song et al., 2025); , and extend to multimodal and vision-conditioned settings (You et al., 2025; Yang et al., 2025), while exploring flexible-length any-order masking (Gong et al., 2025b,a; Song et al., 2025; Yang et al., 2025; You et al., 2025; Kim et al., 2025a). Theoretically, dLLMs can approach the quality of Auto-regressive Models (ARMs) but require multiple denoising steps, where complexity may grow with sequence-level correctness demands and context length (Feng et al., 2025; Liu et al., 2025a). These properties motivate inference acceleration: unlike ARMs, dLLMs lack lossless KV caching and typically need many denoising steps, yielding higher latency (Cobbe et al., 2021; Hendrycks et al., 2021b; Chen et al., 2021; Jain et al., 2024).

**Inference acceleration and decoding strategies.**  To reduce latency, parallel decoding strategies sample multiple tokens per step (Yu et al., 2025; Wei et al., 2025), thereby decreasing the total number of iterations without requiring retraining. In addition, Caching or reusing outputs from bidirectional attention (Yu et al., 2025; Wei et al., 2025) and other stable computations (Liu et al., 2025b) across steps can further reduce latency. Beyond raw speed, planning and ordering improve robustness and sample efficiency: path-planning selects which tokens to (re)mask before denoising; adaptive token ordering learns easier sequences of decisions; PC-Sampler adds position-aware calibration to avoid early trivial selections; DPaD prunes distant suffix attention to curb redundant compute (Peng et al., 2025; Kim et al., 2025b; Huang et al., 2025; Chen et al., 2025).Some methods also take into account the dynamic nature of decoding and leverage historical information to improve decoding performance(Wang et al., 2025). Despite their effectiveness, these score-based methods are largely *history-agnostic*—they rely solely on current-step confidence. This approach can lead to step-level instability and redundant iterations, as tokens are repeatedly re-masked until their confidence converges.

## 3 Preliminary

### 3.1 Inference Process of dLLMs

A Diffusion Large Language Model (dLLM) generates discrete text sequences by iteratively denoising from a fully masked input. Unlike autoregressive (AR) models, which decode tokens sequentially, dLLMs view generation as a stochastic process that starts from a sequence where all positions are masked and gradually recovers the clean sequence $x_0$.
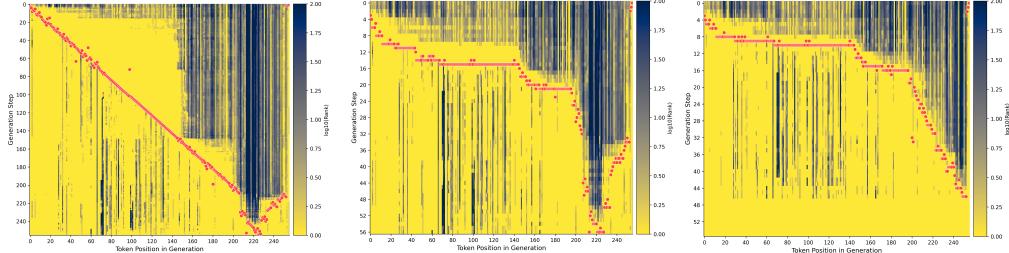
Figure 3: Confidence rank of the **final predicted token** at each position during the generation steps(Left: LLaDA-8B-Inst, Middle: Fast-dLLM, Right: CreditDecoding). The correct token refers to the model's final prediction. Each data point represents the softmax rank of the final output token on a log-scale, color-coded from yellow (top-1) to blue (lower ranks). The red dots denote that the model actually decoded this token at the corresponding step, which is also the *Decoding Boundary* in Figure1.

Let $x_0 = (x_0^1, \ldots, x_0^L) \sim p_{\text{data}}$ be a sequence of length $L$ over a vocabulary $V$ of size $|V|$. At each discrete step $t \in \{0, \ldots, T\}$, denote by $M_t \subseteq \{1, \ldots, L\}$ the set of masked positions. We use $\alpha_t \in [0, 1]$ to denote the proportion of tokens that are *masked* at step $t$, with $\alpha_T = 1$ and $\alpha_0 = 0$. The sequence $\{\alpha_t\}_{t=0}^T$ follows a predefined masking schedule such that $\alpha_t < \alpha_{t+1}$ for all $t$, reflecting the gradual reduction of noise during generation. The core component of a dLLM is a denoising model $f_\theta$, which maps a corrupted sequence $x_t$ to a matrix of logits $f_\theta(x_t) \in \mathbb{R}^{L \times |V|}$. The maximum probability $s_t^i = \max_{v \in V} p_\theta^i(v \mid x_t)$ is referred to as the *confidence score* for position $i$ at step $t$, and the corresponding $\arg\max$ token is the model's current best guess for $x_0^i$.

The training objective is expressed as:

$$\mathcal{L}_{\text{dLLM}} = -\mathbb{E}_{t, x_0, x_t} \left[ \frac{1}{t} \sum_{i=1}^L \mathbf{1}[i \in M_t] \log p_\theta^i(x_0^i \mid x_t) \right] \tag{1}$$

where $p_\theta^i(\cdot \mid x_t) = \text{Softmax}(f_\theta(x_t)^i)$. This order-agnostic formulation enables the model to predict any masked token from arbitrary visible context.

The reverse process starts from the fully masked sequence $x_T$ and iteratively produces less corrupted states until reaching $x_0$. At each step, the update from $x_{t+1}$ to $x_t$ is given by:

$$x_t \sim \prod_{i=1}^L g_\theta(x_t^i \mid x_{t+1}) \tag{2}$$

$$g_\theta(x_t^i \mid x_{t+1}) = \begin{cases} x_{t+1}^i & i \notin M_{t+1}, \\ \text{Cat}\left( \frac{\alpha_t}{\alpha_{t+1}} \mathbf{e}_M + \frac{\alpha_{t+1} - \alpha_t}{\alpha_{t+1}} p_\theta^i(\cdot \mid x_{t+1}) \right) & i \in M_{t+1} \end{cases} \tag{3}$$

Here $\mathbf{e}_M$ denotes the one-hot vector of the mask token. Intuitively, at each masked position, the token either remains masked or be sampled from the model's predictive distribution at the current step. In practice, only a subset of $M_{t+1}$ is selected and decoded at each step (the rest remain masked), iterating until $M_t = \varnothing$.

### 3.2 Parallel Decoding

Unlike autoregressive models, dLLMs can recover multiple positions in parallel at each step, thereby completing sequence generation in fewer iterations and improving efficiency. During inference, predictions at masked positions are often assumed conditionally independent given $x_t$. Thus, at step $t$, for any subset $I \subseteq M_t$, the joint distribution can be approximated by the product of marginals:

$$q_\theta(X_I \mid x_t, t) \approx \prod_{i \in I} p_\theta^i(X^i \mid x_t) \tag{4}$$
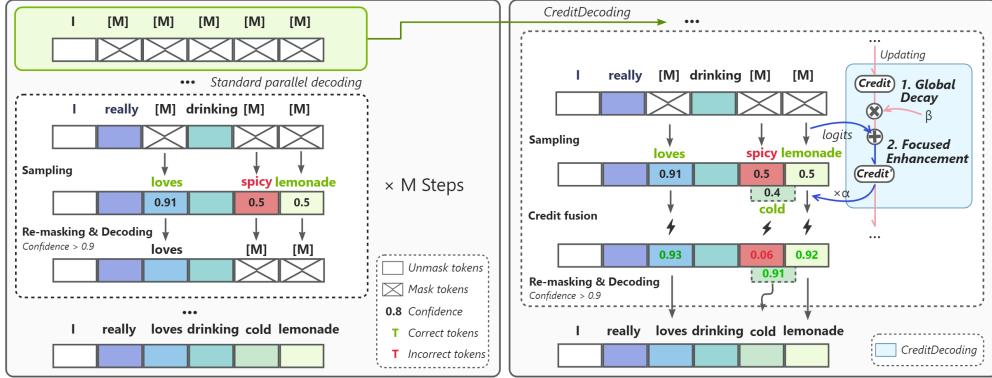
where $X_I = \{X^i\}_{i \in I}$.

4

Figure 4: Comparison between **standard dLLM parallel decoding** (left) and the proposed **Credit-Decoding** (right). The left diagram illustrates how existing methods decode solely on instantaneous predictions at each step, causing the repetitive remasking of correct tokens. In contrast, CreditDecoding maintains a token-level credit value across steps, using Trace Credit as a prior to enhance and calibrate current predictions.

Prior theoretical work (Wu et al., 2025) shows that the product of marginals provides a close approximation to the true joint distribution when the confidence scores $s_t^i$ are sufficiently close to 1, making greedy *parallel* decoding effectively equivalent to *sequential* greedy decoding. In practice, the mainstream implementation of parallel decoding is achieved by setting a confidence threshold and *selecting* positions to update, an approach whose effectiveness has been demonstrated by FastdLLM (Wu et al., 2025):

$$S_t = \{i \in M_t : s_t^i \geq \tau\} \tag{5}$$

where $\tau \in (0, 1)$ is a fixed confidence threshold. Positions $i \in S_t$ are selected to be decoded, while the remaining positions are kept masked.

Beyond maximum probability, other scoring functions have been explored, such as negative entropy or the model's marginal confidence. Likewise, different selection schemes have been proposed: thresholding, top-$k$ selection, or adaptive rules.

## 4 Methodology

### 4.1 Analysis

In this section, we examine the limitations of dLLMs during inference. Figure 3 visualizes, for all tokens that are eventually decoded correctly (i.e., present in the final generated sequence), how their confidence ranks evolve over reverse denoising steps. Many tokens repeatedly appear in the top-1 position long before they are actually decoded. This effect is most pronounced under single-token decoding, which often leads to higher final accuracy (Feng et al., 2025). In this setting, the extremely slow update cycle causes tokens to oscillate between being decoded and re-masked. Even with faster parallel settings, substantial redundant updates remain visible. Figure 6(c)(d) further examines several representative tokens by plotting their confidence trajectories (x-axis: step id; y-axis: confidence; green shading: currently top-1). Correct tokens generally follow a convergent trend, ultimately reaching confidence 1 at the commit step. However, they often linger at low absolute confidence for many steps while already ranking top-1. Threshold-based strategies thus defer commits unnecessarily, incurring extra computation while failing to exploit the signal that these tokens are consistently supported.

Together, these observations highlight two limitations of existing dLLM parallel decoders. **(i) Redundant computation.** Since decoding decisions are gated solely by instantaneous confidence, many correct tokens are repeatedly re-masked until their scores exceed the threshold. This leads to wasted iterations and is particularly severe under single-token decoding, where the process becomes overly conservative. **(ii) History-agnostic decoding.** Each step is treated independently without regard to past predictions. When the model temporarily mispredicts, tokens that were steadily

5

converging may experience confidence fluctuations and fail to be decoded, causing error propagation to subsequent denoising steps.

These limitations motivate our proposed **CreditDecoding**, which accumulates historical model predictions as a token-level prior. By fusing this prior with the current logits, CreditDecoding both accelerates convergence for under-confident but correct tokens and stabilizes decoding against transient fluctuations.

## 4.2 CreditDecoding for Parallel dLLMs Inference

In this section, we introduce *Trace Credit*, a mechanism that tracks the stability of token predictions over time and uses it to estimate the likelihood of convergence to high confidence (ideally close to 1) in subsequent denoising steps. By tracking the historical consistency of model predictions, trace credits serve as a temporal prior that helps distinguish stable, correct hypotheses from transient fluctuations.

**Definition of Trace Credit.** At each reverse diffusion step $t$, for every masked position $i$ and token $v \in \mathcal{V}$, we maintain a credit value $C_t^{i,v} \in \mathbb{R}_{\geq 0}$ that accumulates evidence from past prediction traces. Let $v^* = \arg\max_{v \in \mathcal{V}} p_\theta^i(v \mid x_t)$ denote the current top-1 candidate at position $i$. The credit is updated via a combination of *global decay* and *focused enhancement*:

$$C_t^{i,v} = \begin{cases} \beta\, C_{t+1}^{i,v} + \big(p_\theta^i(v \mid x_t)\big)^\gamma & \text{if } v = v^*, \\ \beta\, C_{t+1}^{i,v} & \text{otherwise,} \end{cases} \qquad \beta \in (0,1),\ \gamma \in (0,1). \tag{6}$$

Here, $\beta$ acts as a decay factor that gradually diminishes older evidence, preventing outdated or incorrect predictions from dominating. The exponent $\gamma < 1$ applies a concave transformation to the current probability, which relatively amplifies low-confidence values, thereby accelerating the stabilization of potentially correct but initially uncertain tokens.

As illustrated in Figure 4, this update rule balances two complementary dynamics. **(i) Global decay:** The discounting factor $\beta$ ensures that stale predictions—especially those with persistently low confidence or tokens not currently favored—are gradually forgotten. This mitigates the risk of error accumulation from spurious, short-lived confidence spikes. **(ii) Focused enhancement:** Only the top-1 predicted token $v^*$ receives an additional credit boost at each step. This focuses the credit on the model's active hypothesis, aligning it with the actual decoding trajectory rather than momentary fluctuations across the full distribution.

Crucially, the accumulated credit is fused with the model's raw logits in the log domain to form an enhanced predictive distribution:

$$\tilde{f}_\theta(x_t)_v^i = f_\theta(x_t)_v^i + \alpha \cdot \log\big(C_t^{i,v} + 1\big), \qquad \alpha > 0, \tag{7}$$

where $f_\theta(x_t)_v^i = \log p_\theta^i(v \mid x_t)$ (up to a constant) denotes the model's original logits, and $\alpha$ controls the strength of the credit-based prior. This fusion is mathematically equivalent to multiplying the original probability by $(1 + C_t^{i,v})^\alpha$, effectively applying a multiplicative prior over the likelihood.
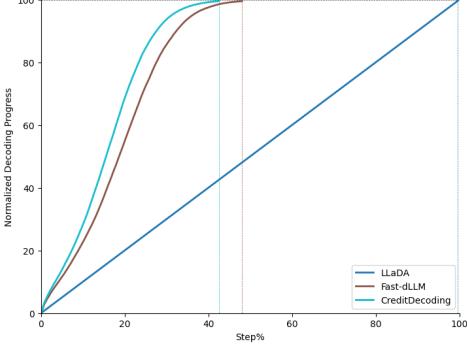
The resulting enhanced distribution is computed as:

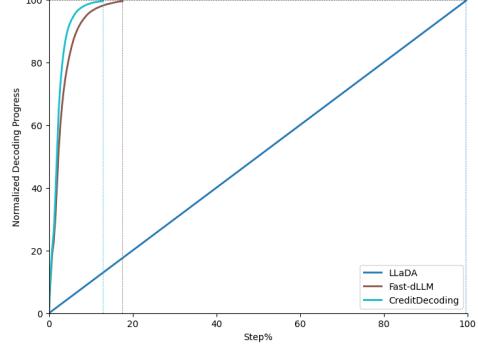$$\tilde{p}^i(v \mid x_t) = \mathrm{Softmax}\Big(\tilde{f}_\theta(x_t)_v^i\Big), \tag{8}$$

which replaces the original $p_\theta^i(v \mid x_t)$ in subsequent sampling and masking decisions. Tokens that have been consistently predicted across steps receive a confidence boost, promoting earlier commitment and reducing redundant recomputation. Conversely, tokens with unstable or sporadic support are suppressed, improving decoding stability—particularly in long-sequence generation and complex reasoning tasks.

Importantly, CreditDecoding does not alter the underlying decoding policy; it merely enhances the model's output distribution. This preserves compatibility with standard inference techniques such as threshold-based sampling, top-$k$ truncation, adaptive KV caching, and compiler-level optimizations, enabling seamless integration and cumulative efficiency gains.

To improve scalability and reduce interference from uncertain future context, we default to maintaining credits only within the current decoding block. This design limits the influence of under-informed positions and enhances robustness across varying sequence lengths and model scales.

(a) Normalized Decoding Progress on GSM8K      (b) Normalized Decoding Progress on SQuAD2.0

Figure 5: Normalized Decoding Progress **w/o Early Stop** on **GSM8K** and **SQuAD2.0**. We demonstrate the *decoding progress* through visualizing the accumulated number of decoded tokens per step, using **LLaDA-8B-Instruct** with *generation length* $= 256$ and *block size* $= 64$.In order to ensure comparability of the decoding progress, we did not use early stopping. The vertical dashed lines in the figure mark the decoding step at which each method completes, and the ratio of these steps represents the speedup.

## 4.3 Full-Distribution Credit Accumulation

The formulation in Equation equation 6 updates credits exclusively for the top-1 candidate at each step. While this focused strategy maximizes signal concentration and accelerates convergence, it may discard potentially useful information from suboptimal but plausible tokens.

We therefore extend CreditDecoding to a more general form that accumulates credits across the entire vocabulary:

$$C_t^{i,v} = \beta \, C_{t+1}^{i,v} + \left( p_\theta^i(v \mid x_t) \right)^\gamma, \quad \forall v \in \mathcal{V}. \tag{9}$$

This maximizes information usage rather than focusing solely on the most likely candidate. Although it may reduce the degree of acceleration (as reinforcement is less concentrated), it improves robustness by retaining more distributional signals. Empirically, this trade-off leads to smaller quality loss while maintaining speedup, making CreditDecoding a flexible framework for balancing efficiency and accuracy.

# 5 Experiments

## 5.1 Experimental Setup

**Implementation** (Nie et al., 2025) and LLaDA-MoE-7B-A1B-Instruct (Zhu et al., 2025b). Specific inference configurations differ across experiments; consequently, our settings may slightly deviate from those reported in the original LLaDA and LLaDA-MoE papers. Nevertheless, we ensured that all other configurations remained consistent before and after integrating CreditDecoding in our experiments.We detail them in the corresponding sections. All experiments are conducted on NVIDIA H20-3e 140 GB GPUs.

In the main experiments, the generation length and number of steps were both set to 256. Additional experiments with different generation lengths are reported in Section 5.5. Based on the analyses in (Appendix C.3 and Section 5.3), we set the block length to 64, and the hyperparameters of Credit Decoding to $\alpha = 0.65$ and $\beta = 0.7$.

**Evaluation Tasks**: In the main experiments, we comprehensively evaluate CreditDecoding on eight datasets spanning five categories. Specifically, we evaluate inference performance on DROP (Dua et al., 2019) and KorBench (Ma et al., 2025), language understanding on SQuAD, knowledge assessment on MMLU (Hendrycks et al., 2021a), coding ability on OpenAI HumanEval (Chen et al., 2021) and LiveCodeBench (Jain et al., 2024), and mathematical reasoning on GSM8K (Cobbe et al., 2021) and Math (Hendrycks et al., 2021b). In the ablation, scaling, and other analysis experiments,

Table 1: Main benchmark results **w/ Early Stop** across eight datasets on LLaDA-8B-Instruct and LLaDA-MoE-Instruct with *a generation length of 256* (CD for Credit Decoding). The generation length and step are fixed at 256, while the block length is set to 64. Each cell presents the performance score (top: score) with the + - values in the bottom-right indicating *the difference relative to LLaDA*, and the tokens per forward (bottom : TPF) with the *relative improvement percentage over Fast-dLLM*.

| Benchmark | LLaDA-8B-Instruct | | | | LLaDA-MoE-Instruct | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | LLaDA | Fast-dLLM | CD | CD* | LLaDA-MoE | Fast-dLLM | CD | CD* |
| MMLU $_{SCORE}$ | 62.46 | $62.43_{-0.03}$ | $\mathbf{63.78}_{+1.32}$ | $63.66_{+1.20}$ | 64.08 | $64.08_{0.00}$ | $\mathbf{64.21}_{+0.13}$ | $63.94_{-0.14}$ |
| $_{TPF}$ | 1 | 2.86 | 4.57 (+56%) | 3.38 (+18%) | 1 | 2.16 | 2.46 (+14%) | 2.33 (+8%) |
| SQuAD2.0 $_{SCORE}$ | 91.43 | $91.43_{0.00}$ | $\mathbf{91.71}_{+0.28}$ | $91.48_{+0.05}$ | 86.88 | $86.88_{0.00}$ | $87.27_{+0.39}$ | $\mathbf{87.35}_{+0.47}$ |
| $_{TPF}$ | 1 | 13.55 | 16.84 (+24%) | 15.07 (+11%) | 1 | 7.09 | 9.64 (+36%) | 8.41 (+19%) |
| DROP $_{SCORE}$ | 82.86 | $82.74_{-0.12}$ | $\mathbf{82.78}_{-0.08}$ | $82.70_{-0.16}$ | **80.16** | $80.16_{0.00}$ | $79.72_{-0.44}$ | $79.87_{-0.29}$ |
| $_{TPF}$ | 1 | 2.93 | 3.79 (+29%) | 3.15 (+8%) | 1 | 2.73 | 3.28 (+20%) | 2.92 (+7%) |
| KorBench $_{SCORE}$ | 33.12 | $33.20_{+0.08}$ | $\mathbf{35.04}_{+1.92}$ | $33.92_{+0.80}$ | 36.72 | $\mathbf{36.88}_{+0.16}$ | $36.48_{-0.24}$ | $36.64_{-0.08}$ |
| $_{TPF}$ | 1 | 3.72 | 5.03 (+35%) | 4.43 (+19%) | 1 | 2.36 | 3.28 (+38%) | 2.73 (+16%) |
| HumanEval $_{SCORE}$ | 34.76 | $34.15_{-0.61}$ | $36.59_{+1.83}$ | $\mathbf{37.80}_{+3.04}$ | 51.22 | $51.22_{0.00}$ | $51.22_{0.00}$ | $\mathbf{53.05}_{+1.83}$ |
| $_{TPF}$ | 1 | 3.82 | 4.69 (+23%) | 4.18 (+9%) | 1 | 4.97 | 6.00 (+21%) | 5.45 (+10%) |
| LCB $_{SCORE}$ | **8.15** | $8.15_{0.00}$ | $7.54_{-0.61}$ | $7.71_{-0.44}$ | 13.88 | $14.04_{+0.16}$ | $14.37_{+0.49}$ | $\mathbf{14.65}_{+0.77}$ |
| $_{TPF}$ | 1 | 1.93 | 2.17 (+12%) | 2.00 (+4%) | 1 | 2.43 | 2.81 (+16%) | 2.55 (+5%) |
| GSM8K $_{SCORE}$ | 77.94 | $\mathbf{78.47}_{+0.53}$ | $77.18_{-0.76}$ | $77.48_{-0.46}$ | 74.37 | $74.45_{+0.08}$ | $\mathbf{74.98}_{+0.61}$ | $74.37_{0.00}$ |
| $_{TPF}$ | 1 | 3.22 | 3.87 (+20%) | 3.39 (+5%) | 1 | 2.28 | 2.68 (+18%) | 2.42 (+6%) |
| Math $_{SCORE}$ | 37.30 | $37.04_{-0.26}$ | $\mathbf{37.24}_{-0.06}$ | $37.18_{-0.12}$ | 36.02 | $35.84_{-0.18}$ | $\mathbf{36.28}_{+0.26}$ | $36.26_{+0.24}$ |
| $_{TPF}$ | 1 | 2.42 | 2.84 (+17%) | 2.55 (+5%) | 1 | 2.35 | 2.71 (+15%) | 2.48 (+6%) |
| **Average** $_{SCORE}$ | 53.50 | $53.45_{-0.05}$ | $53.98_{+0.48}$ | $\mathbf{53.99}_{+0.49}$ | 55.42 | $55.44_{+0.02}$ | $55.57_{+0.15}$ | $\mathbf{55.77}_{+0.35}$ |
| $_{TPF}$ | 1 | 4.31 | 5.48 (+27%) | 4.77 (+11%) | 1 | 3.30 | 4.11 (+25%) | 3.66 (+11%) |

due to computational constraints, we select five representative datasets across categories: MMLU, SQuAD, KorBench, HumanEval, and GSM8K.

*Evaluation metric*: We adopt the standard performance metrics for each evaluation dataset, as detailed in Appendix B. In addition, to examine whether CreditDecoding can mitigate redundant computation inherent in traditional dLLMs, we utilize TPF (Tokens Per Forward) to evaluate dLLM inference efficiency.

**Early Stop**: In dLLM, early stopping changes the generated length and thus significantly affects TPF. Some studies disable it to boost TPF, as the model quickly outputs the EOS token. However, in practical applications it is usually enabled to avoid redundant outputs. Throughout this paper we **enable Early Stop** by default, except in the orthogonal analysis(Figure 2 Left) and Figure 5, where it is disabled for better comparison.

## 5.2 Main Results

As shown in Table 1, we evaluate our methods on eight datasets using LLaDA-8B-Instruct and LLaDA-MoE-Instruct. Here, CD refers to CreditDecoding (Section 4.2), and CD* denotes its extended version (Section 4.3).

The key difference lies in trace credit accumulation: CreditDecoding records only the top-1 logit at each token per step, whereas CreditDecoding* retains all logits to leverage global information from previous steps. We use each benchmark's default performance metric and report TPF for inference speed, with TPF of the baseline normalized to 1. TPF of CreditDecoding and CreditDecoding* thus directly reflects speedup relative to the baseline.

Overall, both CreditDecoding and CreditDecoding* outperform the baseline in performance and speed (Avg column). CreditDecoding achieves a 5.48× speedup and 0.48 performance gain on LLaDA-8B-Instruct, and a 4.10× speedup on LLaDA-MoE-Instruct. CreditDecoding* provides slightly higher performance gains (0.49 and 0.36) with 10% lower speed than CreditDecoding, but still considerably accelerates inference.
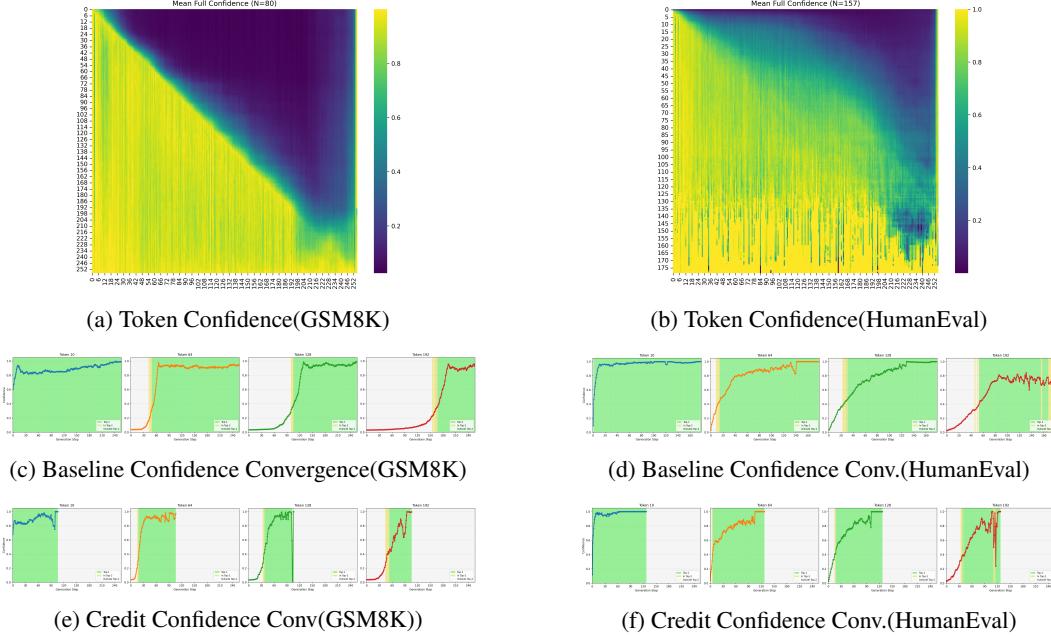
(a) Token Confidence(GSM8K)

(b) Token Confidence(HumanEval)

(c) Baseline Confidence Convergence(GSM8K)

(d) Baseline Confidence Conv.(HumanEval)

(e) Credit Confidence Conv(GSM8K))

(f) Credit Confidence Conv.(HumanEval)

Figure 6: Token confidence across datasets. Figures (a) and (b) present the average confidence of $N$ tokens at each position during inference. Figures (c)–(f) illustrate the convergence behavior of four representative tokens located at positions 10, 64, 128, and 192, respectively. Results in (a), (c), and (e) are sampled from GSM8K, while those in (b), (d), and (f) are from HumanEval. Panels (c) and (d) show baseline convergence, and (e) and (f) show convergence with CreditDecode.

Figure 3 illustrates that after applying CreditDecoding, the red line representing token decoding becomes more horizontal, indicating more parallel decoding within each step. While the baseline requires all 256 steps, CreditDecoding completes decoding in 50 steps, consistent with the $5\times$ speedup in Table 1. Yellow and blue denote high and low confidence, respectively, with their boundary marking the step where the model finalizes predictions. Traditional dLLMs discard remasked token information, necessitating repeated predictions and creating a gap between the red line and the yellow-blue boundary.

Two observations explain CreditDecoding's speedup: it moves the red line closer to the yellow-blue boundary, and the boundary appears earlier, allowing the model to determine decoded tokens more efficiently.

In Figure 5, we visualize the accumulated decoded tokens per step for LLaDA, Fast-dLLM, and CreditDecoding on specific datasets. Two key observations emerge: First, same method speedup varies across datasets. Second, CreditDecoding consistently outperforms Fast-dLLM, especially on SQuAD2.0, where it reduces the decoding steps by an additional 5% step compared to Fast-dLLM's 20% step. As an orthogonal method, CreditDecoding offers greater improvements with higher speedups.

### 5.3 Hyperparameter Ablation Study

As discussed in Section4.2, CreditDecoding has three hyperparameters: $\beta$ for global decay, $\alpha$ for logits fusion, and $\gamma$ for concave amplification. We fix $\gamma$ and perform ablation studies on $\alpha$ and $\beta$ in the range [0, 0.95] with a step size of 0.05.

Figure 7 shows that performance fluctuates slightly with $\alpha$ and $\beta$, peaking around $\beta = 0.5$ and $\beta = 0.7$, while $\alpha$ remains stable within [0.2, 0.65]. Larger values of both parameters increase TPF by accumulating more trace credit, but may reduce accuracy. We select $\alpha = 0.65$ and $\beta = 0.7$ as they provide a good trade-off, yielding strong performance and high TPF across five datasets, though optimal values vary by dataset.
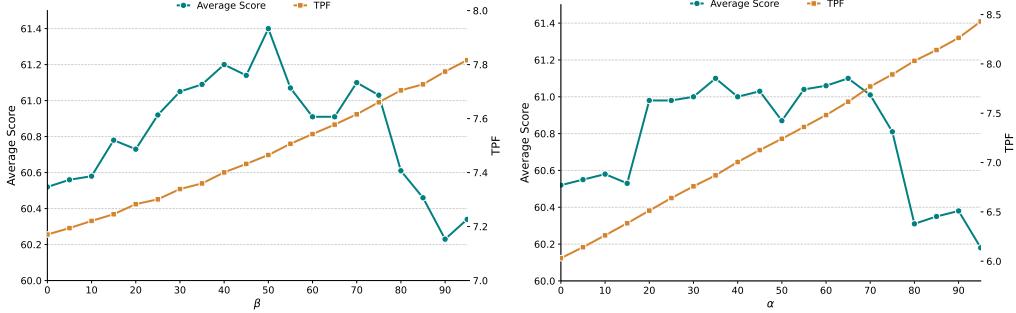
Figure 7: Hyperparameter ablation. The blue lines in the left and right figures show the ablation of $\beta$ and $\alpha$ from 0 to 0.95 with a step size of 0.05. In the left figure, $\alpha$ is fixed at 0.65, and in the right figure, $\beta$ is set to 0.7. The green line represents the performance of LLaDA-8B-Instruct, and the red line indicates the TPF under the corresponding hyperparameter settings.

Further analysis in Figure 6 compares GSM8K and HumanEval. GSM8K requires more context from prior steps, leading to delayed but sharp confidence gains, while HumanEval shows earlier, gradual confidence increases. With CreditDecoding, HumanEval retains its confidence trend while reducing inference steps, whereas GSM8K experiences minor fluctuations that contribute to the performance drop reported in Table 1.

## 5.4 Scalability

Current research on dLLMs typically evaluates generation lengths of 128 or 256, with a few studies extending to 512. However, the primary advantage of dLLMs over autoregressive models—parallel decoding—is largely underutilized at such short lengths. To address this, in this section we scale the generation length to 1024 and 4096, aiming to provide insights into the potential of dLLMs for long-text generation.

We fix the number of steps equal to the generation length and set the block length to 64. As shown in Figure 2, both the baseline and CreditDecoding achieve their best scores at length 512, with performance gradually degrading as the length increases. Notably, CreditDecoding consistently outperforms the baseline, and its performance declines more slowly with increasing generation length, demonstrating stronger robustness for long-text generation.

In terms of inference speed, constrained by the block length, the TPF speedup remains around $7\times$. However, as the total inference time increases with longer texts, the practical acceleration benefits of CreditDecoding become even more pronounced.

## 5.5 Orthogonality

CreditDecoding operates purely on the model logits and does not interfere with the *sample* or *select* procedures. This design makes it a plug-and-play post-processing module that is naturally orthogonal to both (i) system-level inference optimizations such as compiler-level acceleration ((Jain et al., 2023)), and (ii) algorithmic acceleration strategies for dLLMs such as threshold decoding, KV-cache variants, and EOS early stopping.

In all cases, we keep the hyperparameters and selection rules of the baseline methods unchanged, and only apply CreditDecoding on top. As demonstrated in Figure 2, our experiments confirm this orthogonality. When combined with CreditDecoding, we observe that the speedup is preserved while the performance drop is partially mitigated, showing that historical trace credit acts as a stabilizing prior under aggressive compression.

For algorithmic acceleration, threshold parallel decoding (Fast-dLLM w/o KV cache (Wu et al., 2025)) typically accelerates but at the cost of lower accuracy. For EOS early stopping, CreditDecoding complements the strategy by improving robustness of token selection in earlier steps. Adding CreditDecoding further boosts both speed and accuracy, with $4.7\times$ acceleration over the baseline with an additional $+0.63$ improvement in average score. Similarly, CreditDecoding alleviates the accuracy

loss of threshold decoding and complements KV-cache methods by reducing redundant iterations within cached segments. These results demonstrate that CreditDecoding can be seamlessly integrated with existing optimizations, offering consistent gains without modifying sampling or selection logic.

Details of the acceleration methods we used are provided in Appendix C.4, along with orthogonality experiment results for CreditDecoding on LLaDA-8B-Instruct and LLaDA-MoE-Instruct, including TPS inference speed. Additionally, orthogonality experiments for FP8 quantization on LLaDA-MoE-Instruct are also presented.

## 6 Conclusion

The history-agnostic decoding process of diffusion language models introduces substantial computational redundancy. We propose CreditDecoding, a training-free decoding strategy orthogonal to mainstream acceleration methods. For each token to be decoded, CreditDecoding assigns a trace credit score based on the model's logits predictions from historical steps and applies this correction to the current logits.

By fully leveraging the model's past predictions on remasked tokens for the first time, CreditDecoding reduces the complexity of current inference steps, leading to improvements in both performance and decoding efficiency, resulting in a $5.48\times$ speedup and an average accuracy gain of $+0.48$ compared to LLaDA-8B-Instruct.

The goal of CreditDecoding is to approach the theoretical limit of acceleration(the yellow-blue boundary in Figures6(a) and (b)). As more powerful base models emerge, the acceleration ceiling of CreditDecoding will increase, leading to further efficiency gains in model inference.

# References

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021. URL https://arxiv.org/abs/2107.03374.

Xinhua Chen, Sitao Huang, Cong Guo, Chiyue Wei, Yintao He, Jianyi Zhang, Hai "Helen" Li, and Yiran Chen. Dpad: Efficient diffusion language models with suffix dropout, 2025. URL https://arxiv.org/abs/2508.14148.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems, 2021. URL https://arxiv.org/abs/2110.14168.

Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. Drop: A reading comprehension benchmark requiring discrete reasoning over paragraphs, 2019. URL https://arxiv.org/abs/1903.00161.

Guhao Feng, Yihan Geng, Jian Guan, Wei Wu, Liwei Wang, and Di He. Theoretical benefit and limitation of diffusion language model, 2025. URL https://arxiv.org/abs/2502.09622.

Shansan Gong, Shivam Agarwal, Yizhe Zhang, Jiacheng Ye, Lin Zheng, Mukai Li, Chenxin An, Peilin Zhao, Wei Bi, Jiawei Han, Hao Peng, and Lingpeng Kong. Scaling diffusion language models via adaptation from autoregressive models, 2025a. URL https://arxiv.org/abs/2410.17891.

Shansan Gong, Ruixiang Zhang, Huangjie Zheng, Jiatao Gu, Navdeep Jaitly, Lingpeng Kong, and Yizhe Zhang. Diffucoder: Understanding and improving masked diffusion models for code generation, 2025b. URL https://arxiv.org/abs/2506.20639.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding, 2021a. URL https://arxiv.org/abs/2009.03300.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset, 2021b. URL https://arxiv.org/abs/2103.03874.

Pengcheng Huang, Shuhao Liu, Zhenghao Liu, Yukun Yan, Shuo Wang, Zulong Chen, and Tong Xiao. Pc-sampler: Position-aware calibration of decoding bias in masked diffusion models, 2025. URL https://arxiv.org/abs/2508.13021.

Animesh Jain, Shunting Zhang, Edward Yang, et al. Pytorch 2.0: Our next generation 2.0 release, 2023. URL https://arxiv.org/abs/2305.01916.

Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code, 2024. URL https://arxiv.org/abs/2403.07974.

Jaeyeon Kim, Lee Cheuk-Kit, Carles Domingo-Enrich, Yilun Du, Sham Kakade, Timothy Ngotiaoco, Sitan Chen, and Michael Albergo. Any-order flexible length masked diffusion, 2025a. URL https://arxiv.org/abs/2509.01025.

Jaeyeon Kim, Kulin Shah, Vasilis Kontonis, Sham Kakade, and Sitan Chen. Train for the worst, plan for the best: Understanding token ordering in masked diffusions, 2025b. URL `https://arxiv.org/abs/2502.06768`.

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention, 2023. URL `https://arxiv.org/abs/2309.06180`.

Xiaoran Liu, Zhigeng Liu, Zengfeng Huang, Qipeng Guo, Ziwei He, and Xipeng Qiu. Longllada: Unlocking long context capabilities in diffusion llms, 2025a. URL `https://arxiv.org/abs/2506.14429`.

Zhiyuan Liu, Yicun Yang, Yaojie Zhang, Junjie Chen, Chang Zou, Qingyuan Wei, Shaobo Wang, and Linfeng Zhang. dllm-cache: Accelerating diffusion large language models with adaptive caching, 2025b. URL `https://arxiv.org/abs/2506.06295`.

Kaijing Ma, Xinrun Du, Yunran Wang, Haoran Zhang, Zhoufutu Wen, Xingwei Qu, Jian Yang, Jiaheng Liu, Minghao Liu, Xiang Yue, Wenhao Huang, and Ge Zhang. Kor-bench: Benchmarking language models on knowledge-orthogonal reasoning tasks, 2025. URL `https://arxiv.org/abs/2410.06526`.

Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-Rong Wen, and Chongxuan Li. Large language diffusion models, 2025.

Fred Zhangzhi Peng, Zachary Bezemek, Sawan Patel, Jarrid Rector-Brooks, Sherwood Yao, Avishek Joey Bose, Alexander Tong, and Pranam Chatterjee. Path planning for masked diffusion model sampling, 2025. URL `https://arxiv.org/abs/2502.03540`.

Yuxuan Song, Zheng Zhang, Cheng Luo, Pengyang Gao, Fan Xia, Hao Luo, Zheng Li, Yuehang Yang, Hongli Yu, Xingwei Qu, Yuwei Fu, Jing Su, Ge Zhang, Wenhao Huang, Mingxuan Wang, Lin Yan, Xiaoying Jia, Jingjing Liu, Wei-Ying Ma, Ya-Qin Zhang, Yonghui Wu, and Hao Zhou. Seed diffusion: A large-scale diffusion language model with high-speed inference, 2025. URL `https://arxiv.org/abs/2508.02193`.

Wen Wang, Bozhen Fang, Chenchen Jing, Yongliang Shen, Yangyi Shen, Qiuyu Wang, Hao Ouyang, Hao Chen, and Chunhua Shen. Time is a feature: Exploiting temporal dynamics in diffusion language models, 2025. URL `https://arxiv.org/abs/2508.09138`.

Qingyan Wei, Yaojie Zhang, Zhiyuan Liu, Dongrui Liu, and Linfeng Zhang. Accelerating diffusion large language models with slowfast sampling: The three golden principles, 2025. URL `https://arxiv.org/abs/2506.10848`.

Chengyue Wu, Hao Zhang, Shuchen Xue, Zhijian Liu, Shizhe Diao, Ligeng Zhu, Ping Luo, Song Han, and Enze Xie. Fast-dllm: Training-free acceleration of diffusion llm by enabling kv cache and parallel decoding, 2025. URL `https://arxiv.org/abs/2505.22618`.

Ling Yang, Ye Tian, Bowen Li, Xinchen Zhang, Ke Shen, Yunhai Tong, and Mengdi Wang. Mmada: Multimodal large diffusion language models, 2025. URL `https://arxiv.org/abs/2505.15809`.

Jiacheng Ye, Zhihui Xie, Lin Zheng, Jiahui Gao, Zirui Wu, Xin Jiang, Zhenguo Li, and Lingpeng Kong. Dream 7b: Diffusion large language models, 2025. URL `https://arxiv.org/pdf/2508.15487v1`.

Zebin You, Shen Nie, Xiaolu Zhang, Jun Hu, Jun Zhou, Zhiwu Lu, Ji-Rong Wen, and Chongxuan Li. Llada-v: Large language diffusion models with visual instruction tuning, 2025. URL `https://arxiv.org/abs/2505.16933`.

Runpeng Yu, Xinyin Ma, and Xinchao Wang. Dimple: Discrete diffusion multimodal large language model with parallel decoding, 2025. URL `https://arxiv.org/abs/2505.16990`.

Fengqi Zhu, Rongzhen Wang, Shen Nie, Xiaolu Zhang, Chunwei Wu, Jun Hu, Jun Zhou, Jianfei Chen, Yankai Lin, Ji-Rong Wen, et al. Llada 1.5: Variance-reduced preference optimization for large language diffusion models, 2025a.

Fengqi Zhu, Zebin You, Yipeng Xing, Zenan Huang, Lin Liu, Yihong Zhuang, Guoshan Lu, Kangyu Wang, Xudong Wang, Lanning Wei, Hongrui Guo, Jiaqi Hu, Wentao Ye, Tieyuan Chen, Chenchen Li, Chengfu Tang, Haibo Feng, Jun Hu, Jun Zhou, Xiaolu Zhang, Zhenzhong Lan, Junbo Zhao, Da Zheng, Chongxuan Li, Jianguo Li, and Ji-Rong Wen. Llada-moe: A sparse moe diffusion language model, 2025b. URL `https://arxiv.org/abs/2509.24389`.

# A Reproducibility

We employed two widely adopted and advanced dLLM models, LLaDA-8B-Instruct and LLaDA-MoE-Instruct, along with their publicly available weights. Detailed experimental configurations for each setup are provided, and the specific metric configurations for score evaluation are included in Appendix B. We utilized TPF, a metric that is relatively robust to confounding factors, as it is not influenced by the number of GPUs, hardware model, or inference framework. Therefore, we argue that the results reported in this paper are highly reproducible.

# B Evaluation Config

We employed OpenCompass to assist in the evaluation process, ensuring a standardized and systematic assessment. For each benchmark(LCB for LiveCodeBench), we utilized the specific metrics presented in Table 2, which allowed for a consistent and comprehensive comparison across different tasks. Regarding in-context learning (ICL) configurations, all benchmarks were evaluated in a zero-shot setting, except for Drop and SQuAD2, which were evaluated in two-shot and one-shot settings respectively.

Table 2: Benchmarks, their corresponding evaluation metrics, and In-Context Learning (ICL) setup.

| Benchmark | Metric | ICL |
|---|---|---|
| GSM8K | Accuracy | 0-shot |
| Math | Accuracy | 0-shot |
| SQuAD2.0 | Score | 1-shot |
| DROP | Score | 2-shot |
| MMLU | Weighted Average | 0-shot |
| KorBench | Naive Average | 0-shot |
| LCB OC Code Generation v6 | Score | 0-shot |
| OpenAI HumanEval | Pass@1 | 0-shot |

# C CreditDecoding Supplement

In this section, we provide additional details and experiments related to CreditDecoding that were not presented in the main body of the paper.

## C.1 Algorithm

To accurately illustrate the workflow of the CreditDecoding algorithm, we present its detailed steps in Algorithm 1. Here, $T$ denotes the decoding step, and $r,s,t$ correspond to the mask ratios of the previous, current, and subsequent steps, respectively. Lines 4–12 describe the process of trace credit accumulation, where historical trace credits are decayed by a factor $\beta$ and the token with the highest confidence is assigned additional trace credit. Lines 13–14 show the application of trace credit, in which past information is integrated into the current step by fusing it with the original logits. Finally, Lines 16–23 represent the mainstream approach to parallel decoding, namely the thresholding method, whose effectiveness has been demonstrated in Fast-dLLM(Wu et al., 2025).

**Algorithm 1** CreditDecode

---

1: Initialize $\mathbf{x}_T \leftarrow ([\text{MASK}], \ldots, [\text{MASK}]), \quad \mathbf{C}_T \leftarrow \mathbf{0}$     ▷ Start from full masking with zero credits
2: **for** $k = T-1, T-2, \ldots, 0$ **do**     ▷ Reverse diffusion step $k$
                                                             ▷ **Credit Update (Eq. equation 6)**
3:     **for** each masked position $i$ **do**
4:         $v^* \leftarrow \arg\max_v p_\theta(v \mid x_{k+1})$     ▷ Top-1 candidate
5:         **for** each $v \in \mathcal{V}$ **do**
6:             $C_k^{i,v} \leftarrow \beta \cdot C_{k+1}^{i,v}$     ▷ Global decay
7:             **if** $v = v^*$ **then**
8:                 $C_k^{i,v} \leftarrow C_k^{i,v} + \left(p_\theta(v \mid x_{k+1})\right)^\gamma$     ▷ Focused enhancement
9:             **end if**
10:         **end for**
11:     **end for**
                                                             ▷ **Logits Fusion (Eq. equation 7)**
12:     **for** each masked position $i$ **do**
13:         **for** each $v \in \mathcal{V}$ **do**
14:             $\tilde{f}_\theta^i(v) \leftarrow f_\theta^i(v) + \alpha \cdot \log(1 + C_k^{i,v})$
15:         **end for**
16:         $\tilde{p}^i(\cdot) \leftarrow \text{Softmax}\left(\tilde{f}_\theta^i(\cdot)\right)$
17:     **end for**
                                                             ▷ **Decoding Decision**
18:     **for** each masked position $i$ **do**
19:         **if** $\max_v \tilde{p}^i(v) > \lambda$ **then**
20:             $x_k^i \sim \text{Cat}\left(\tilde{p}^i(\cdot)\right)$     ▷ Decode tokens
21:         **else**
22:             $x_k^i \leftarrow [\text{MASK}]$     ▷ Re-masked tokens
23:         **end if**
24:     **end for**
25: **end for**
26: **return** $\mathbf{x}_0$: the final generated sequence

---

## C.2 Ideal Decoding Boundary

In Figure 1, we illustrate the decoding boundaries of three methods—LLaDA, Fast-dLLM, and CreditDecoding-within a single plot. Figure 3 further presents the decoding boundaries (highlighted by red lines) from the perspective of confidence rank, along with an analysis of the ideal decoding boundary. For a more intuitive comparison, the ranges of the plots in Figure 1 are aligned, whereas the vertical scales of the three subplots in Figure 3 differ to emphasize the detailed variations introduced by thresholding. The line connecting the first appearance of yellow points, which indicate the initial top-1 confidence, represents the ideal decoding boundary.

It is worth noting that, as shown in Figure 3, different methods yield different ideal decoding boundaries. This variation arises because each method receives distinct inputs at every step, resulting in different confidence distributions. Importantly, stronger models or more effective methods tend to produce an ideal decoding boundary that shifts upward. Our CreditDecoding approach, however, is orthogonal to most existing methods and can further improve their ideal decoding boundaries, bringing the actual decoding boundary closer to the ideal one.

## C.3 Block Length Ablation

In the parallel decoding framework of generative models, the choice of generation block length directly determines the balance between system performance and efficiency. Through block length ablation experiments analyzing Fast-dLLM and CreditDecoding, as shown in Figure 8, it can be observed that when the block length is set to 64, both methods demonstrate relatively superior comprehensive performance. Under this configuration, both CreditDecoding and Fast-dLLM achieve excellent average scores, indicating that this scale fully unleashes the potential of block-level generation mechanisms without significantly sacrificing generation quality. In contrast, smaller block lengths

(such as 32) can maintain high performance levels but result in significantly lower decoding speeds, making it difficult to meet the high-throughput requirements of practical deployment. Meanwhile, larger block lengths (such as 256 and above) significantly improve TPF but lead to a sharp performance decline. Therefore, a block length of 64 achieves the optimal balance between maintaining generation performance and achieving parallel acceleration, making it an ideal choice that considers both accuracy and practicality.

Further comparison between CreditDecoding and Fast-dLLM, as illustrated in Figure 8, reveals that CreditDecoding demonstrates greater advantages in overall performance. In terms of performance, CreditDecoding consistently performs on par with or slightly better than Fast-dLLM under medium and small block lengths , reaching its performance peak at a block length of 32, which demonstrates stronger capabilities. Even when both methods experience performance degradation under large block lengths, CreditDecoding still exhibits better error control. Regarding speed, CreditDecoding shows higher TPF across all block length scales, and the performance gap widens as block length increases, reflecting its architectural advantages in parallel scenarios. In summary, CreditDecoding's core value lies in its superior generation performance compared to Fast-dLLM, along with higher overall efficiency. Future work could explore adaptive block sizing or hybrid decoding strategies to dynamically adjust block length and further optimize the balance between performance and speed across different sequence lengths.
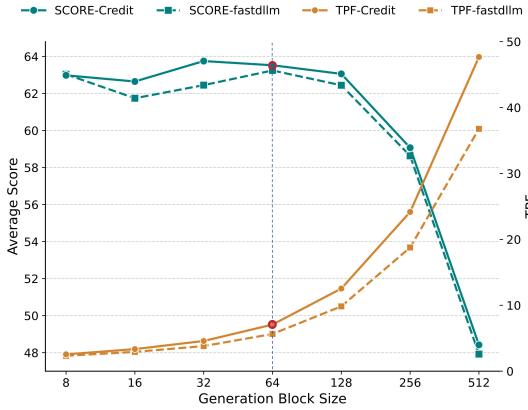


Figure 8: Abaltion Study on Block Length

## C.4 Orthogonality

In Section 5.5, we demonstrate the orthogonality and compatibility of CreditDecoding through experiments combining it with several acceleration techniques. Results show that CreditDecoding consistently improves both speed and performance across all tested methods.

In this section, we provide brief introductions to the acceleration methods discussed in Section 5.5 and include additional TPS results to better illustrate CreditDecoding's effectiveness, particularly on system-level accelerations that mainly improve TPS. We also extend our orthogonality analysis to LLaDA-MoE, with detailed results presented below.

We evaluate its orthogonality on four representative acceleration techniques, as illustrated in Figure 2.

**Early Stop**: Early Stop terminates decoding when the current token is <EOS> and all previous tokens are finalized, effectively reducing redundant generation and improving decoding efficiency.

**Fast-dLLM** (Wu et al., 2025): A state-of-the-art acceleration method consisting of threshold-based parallel decoding and KV Cache. Since the KV Cache significantly increases TPS at the cost of performance, we mainly compare with Fast-dLLM (w/o KV).

**PyTorch Compiler** (Jain et al., 2023): PyTorch Compiler leverages graph-level optimizations for runtime acceleration without altering decoding behavior.

**FP8 Quantization** (Kwon et al., 2023): FP8 quantization is a technique that reduces the precision of floating-point numbers to 8-bit, aiming to accelerate deep learning models by lowering storage and computation costs while maintaining sufficient accuracy.

Table 3: Orthogonality Experiment **w/o Early Stop** on LLaDA-8B-Instruct.

| Method | TPS | TPF | Score |
|---|---|---|---|
| LLaDA-8B-Ins$_{wCD}^{w/oCD}$ | 7.95 | 1.00 | 60.12 |
| | 34.90$_{+339\%}$ | 15.39$_{+1439\%}$ | 60.73$_{+0.61}$ |
| Fast-dLLM (w/o KV)$_{wCD}^{w/oCD}$ | 28.90 | 12.64 | 60.11 |
| | 34.90$_{+21\%}$ | 15.39$_{+22\%}$ | 60.73$_{+0.62}$ |
| Fast-dLLM (w KV)$_{wCD}^{w/oCD}$ | 39.38 | 4.42 | 58.51 |
| | 51.40$_{+31\%}$ | 14.00$_{+217\%}$ | 58.63$_{+0.12}$ |
| Early Stop$_{wCD}^{w/oCD}$ | 9.70 | 1.00 | 59.94 |
| | 37.33$_{+285\%}$ | 6.98$_{+598\%}$ | 60.75$_{+0.81}$ |
| Pytorch Compiler$_{wCD}^{w/oCD}$ | 9.03 | 1.00 | 60.26 |
| | 39.51$_{+337\%}$ | 15.41$_{+1441\%}$ | 60.43$_{+0.17}$ |

Table 4: Orthogonality Experiment **w/o Early Stop** on LLaDA-MoE-Instruct.

| Method | TPS | TPF | Score |
|---|---|---|---|
| LLaDA-MoE-Ins$_{wCD}^{w/oCD}$ | 3.53 | 1.00 | 62.73 |
| | 15.26$_{+333\%}$ | 14.80$_{+1380\%}$ | 62.97$_{+0.24}$ |
| Fast-dLLM (w/o KV)$_{wCD}^{w/oCD}$ | 13.30 | 13.08 | 62.74 |
| | 15.26$_{+15\%}$ | 14.80$_{+13\%}$ | 62.97$_{+0.23}$ |
| FP8 Quantization$_{wCD}^{w/oCD}$ | 2.72 | 1.00 | 62.47 |
| | 11.87$_{+336\%}$ | 14.45$_{+1345\%}$ | 62.58$_{+0.11}$ |
| Early Stop$_{wCD}^{w/oCD}$ | 5.11 | 1.00 | 62.66 |
| | 16.99$_{+233\%}$ | 4.77$_{+377\%}$ | 62.92$_{+0.26}$ |
| Pytorch Compiler$_{wCD}^{w/oCD}$ | 2.42 | 1.00 | 63.29 |
| | 10.72$_{+343\%}$ | 14.88$_{+1388\%}$ | 62.99$_{-0.3}$ |

Table 3 presents detailed results corresponding to Figure 2, including TPS comparisons. Table 4 reports results under the same settings on LLaDA-MoE, further including FP8 quantization.