

Diffusion LLM with Native Variable Generation Lengths: Let [EOS] Lead the Way

Yicun Yang¹ * Cong Wang¹ Shaobo Wang¹ Zichen Wen¹ Biqing Qi²
Hanlin Xu³ Linfeng Zhang¹ †

Shanghai Jiao Tong University¹ Shanghai AI Lab² Huawei³

Abstract

Diffusion-based large language models (dLLMs) have exhibited substantial potential for parallel text generation, which may enable more efficient generation compared to autoregressive models. However, current dLLMs suffer from **fixed generation lengths**, which indicates the generation lengths of dLLMs have to be determined before decoding as a hyper-parameter, leading to issues in efficiency and flexibility. To solve these problems, in this work, we propose to train a diffusion LLM with native variable generation lengths, abbreviated as **dLLM-Var**. Concretely, we aim to train a model to accurately predict the [EOS] token in the generated text, which makes a dLLM be able to natively infer in a block diffusion manner, while still maintaining the ability of global bi-directional (full) attention and high parallelism. Experiments on standard benchmarks demonstrate that our method achieves a **30.1×** speedup over traditional dLLM inference paradigms and a **2.4×** speedup relative to autoregressive models such as Qwen and Llama. Our method achieves higher accuracy and faster inference, elevating dLLMs beyond mere academic novelty and supporting their practical use in real-world applications. *Codes and models have been released.*

🤗 Huggingface model: <https://huggingface.co/maomaocun/dLLM-Var>
🐙 Github repo: <https://github.com/maomaocun/dLLM-Var>



Figure 1 | **Overview of Probabilistic Modeling Paradigms for Text Generation:** Evolution from Autoregressive to Diffusion-Based Approaches. AR & MTP (left): Low parallelism, variable generation. **Vanilla dLLMs** (right): High parallelism, fixed lengths. **dLLM-Var** (middle): variable generation lengths while maintaining parallelism.

1. Introduction

Recently, a multitude of diffusion-based large language models (dLLMs) have emerged [1–9], positioning themselves as strong competitors to GPT-like models in the realm of scaling up. Concurrently, the pursuit of efficient and straightforward inference mechanisms for dLLMs represents a highly promising research avenue. Compared with the multi-token prediction [10] and speculative decoding in autoregressive models, dLLMs have demonstrated the native ability of parallel decoding without any additional modules. Recent works including

¹Project head: yangyicun187@gmail.com

²Corresponding author: zhanglinfeng@sjtu.edu.cn

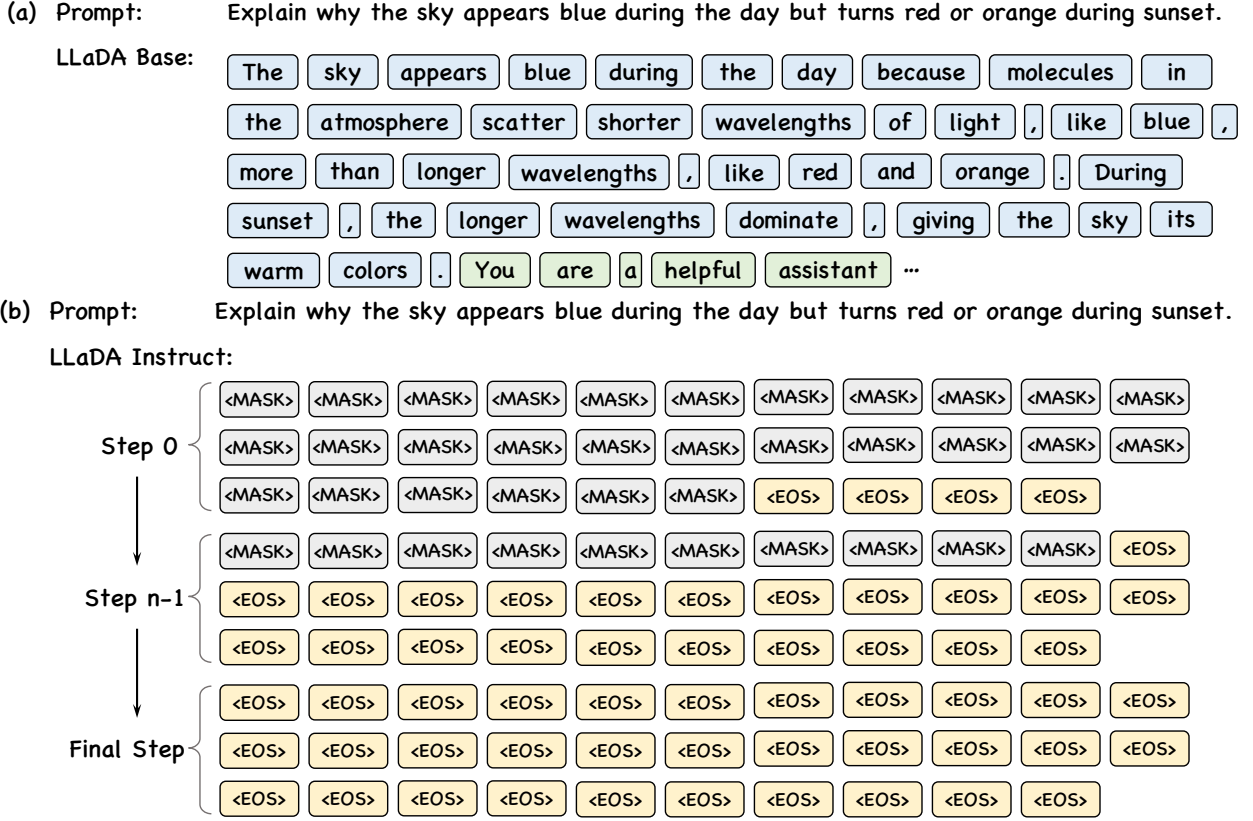


Figure 2 | (a) LLaDA Base: generates irrelevant content without timely EOS token, complicating response extraction. (b) LLaDA Instruct under pure diffusion: fails to produce any effective text, prematurely filling trailing masks with EOS tokens.

SDAR [11], Fast-dLLM v2 [9], and D2F [12] have shown significantly higher decoding speeds than AR models. However, success comes with still great challenges:

Diffusion LLMs suffer from Fixed Length: dLLMs suffer from fixed generation lengths, which indicates that the generation length of them must be pre-defined before decoding as a hyper-parameter, which leads to abundant problems. For instance, an overlong generation length usually leads to repetition and hallucination, as well as abundant computation costs. On the other hand, an overly short generation length tends to have a negative influence on model accuracy, especially in reasoning tasks. Besides, in some applications, such as OCR, it is impossible to estimate the generation lengths before computation. As a result, the fixed generation lengths have limited the practical application of dLLMs.

Block Diffusion Models Suffer from Lower Parallelism. To solve this problem, several recent works have introduced the block diffusion models, where tokens in the same block are parallel decoded while blocks are organized in an autoregressive manner. Block diffusion models achieve the ability of any-length generation by inheriting the block-wise autoregressive generation. However, their parallel potential is severely constrained by their block sizes, which must be decided during training. Besides, as discussed by SDAR, it can be difficult to train a block diffusion with a very large block size, which further limits its upper bounds of parallelism. Besides, since the tokens in the following blocks do not influence tokens in previous works, it is not possible for block diffusion models to perform self-correction, which is a valuable ability for dLLMs.

The previous problems raise the question: *is it possible to train a dLLM with variable generation length while still maintaining great parallelism?* Ideally, this can be achieved by inferring a dLLM in a block diffusion manner, which indicates continuously padding a block of <mask> tokens until the generation of the [EOS] token (*i.e.*, the special token for termination), while still maintaining the bidirectional attention of dLLM. To achieve this, the model should be capable of accurately identifying the [EOS] token during decoding. In this paper, we first dive into the failure of [EOS] generation for the base model and the instruction model of LLaDA. As shown in

Figure 2, the base model usually can not generate the [EOS] token even if it has completed its answer. Instead, it usually keeps generating grammatically correct tokens but makes no sense. For the model after SFT (*i.e.*, the instruction model), without special limitation on the decoding tokens, the model usually tends to first decode all tokens in the tail to [EOS] tokens and then gradually generate [EOS] tokens at the beginning positions, losing the ability to generate any meaningful words. These phenomena demonstrate that current dLLMs do not have the ability for precise [EOS] prediction, which further explains their limitation on fixed-length generation.

To address this problem, this paper proposes to train a dLLM with native ability of variable-length generation (dLLM-Var) by making it accurately decode the [EOS] token. Building on this insight, we introduce a novel training paradigm that empowers dLLMs to replicate the generative efficacy of block diffusion [13] under full-attention regimes, while harnessing key-value (KV) caching for inference. dLLM-Var consists of two strategies. Firstly, we employ a deterministic noise scheduling for [EOS] tokens, where the EOS token is consistently replaced with a [MASK] token during training. Secondly, we propose to package multiple samples into a single sequence and train the dLLM over them without special attention masks, which enables the dLLM to understand the contextual function of the termination token, making it be able to predict [EOS] token at the proper position *i.e.*, the position that a context has been ended while a unrelated context will begin.

In summary, this paper has the following contributions.

- We propose dLLM-Var, which introduce two SFT training strategies for dLLM, making a pretrained dLLM able to be inferred in a block diffusion manner, achieving native variable-length generation. Compared with the traditional block diffusion model, it maintains the bi-directional attention, making it possible for more flexible applications such as editing.
- Based on the block diffusion-style generation, dLLM becomes more compatible with KV cache, avoiding the requirements for complex KV cache design, such as KV refreshing, delaying, and recomputing in previous works, leading to better efficiency and more elegant implementation.
- Extensive experiment results demonstrate the performance of dLLM-Var, achieving comparable and even better performance than a dLLM with optimal fixed-generation length with over $30.1\times$ acceleration, $2.4\times$ faster than AR models. Its efficiency can be further improved based on step distillation [14].

2. Related Work

2.1. Diffusion-based Large Language Models

Diffusion-based large language models (dLLMs) have recently gained prominence as scalable alternatives to autoregressive transformers, leveraging iterative denoising processes for text generation [1–6]. These models offer advantages in parallelism but face challenges in handling variable-length sequences and efficient caching mechanisms. Early works focused on adapting diffusion principles from continuous domains to discrete text spaces, enabling non-autoregressive generation with reduced latency compared to GPT-like architectures.

2.2. KV Cache Management in dLLMs

Key-value (KV) cache optimization is critical for efficient inference in large language models, and dLLMs introduce unique challenges due to their iterative sampling nature. Techniques such as dLLM-Cache [15] explore KV cache refreshing to mitigate redundancy in diffusion steps, while Fast-dLLM [11] introduces training-free acceleration via block-wise updates. Similarly, dKV-Cache [16] proposes dual caching strategies to balance memory usage and generation quality. These methods enhance KV cache utilization in dLLMs [17, 18], yet they often require custom operators and struggle with fixed-length constraints.

2.3. Flexible-Length Generation

Facilitating variable-length outputs is crucial for the practical deployment of generative models. Methods such as DreamOn [19], DAEDAL [20], and FlexMDMs [21] employ specialized tokens to dynamically elongate sequences. However, these techniques introduce overheads arising from diminished effective token density and protracted inference times. Such token-centric extensions underscore the inherent trade-offs between preserving diffusion stability and enabling unbounded generation.

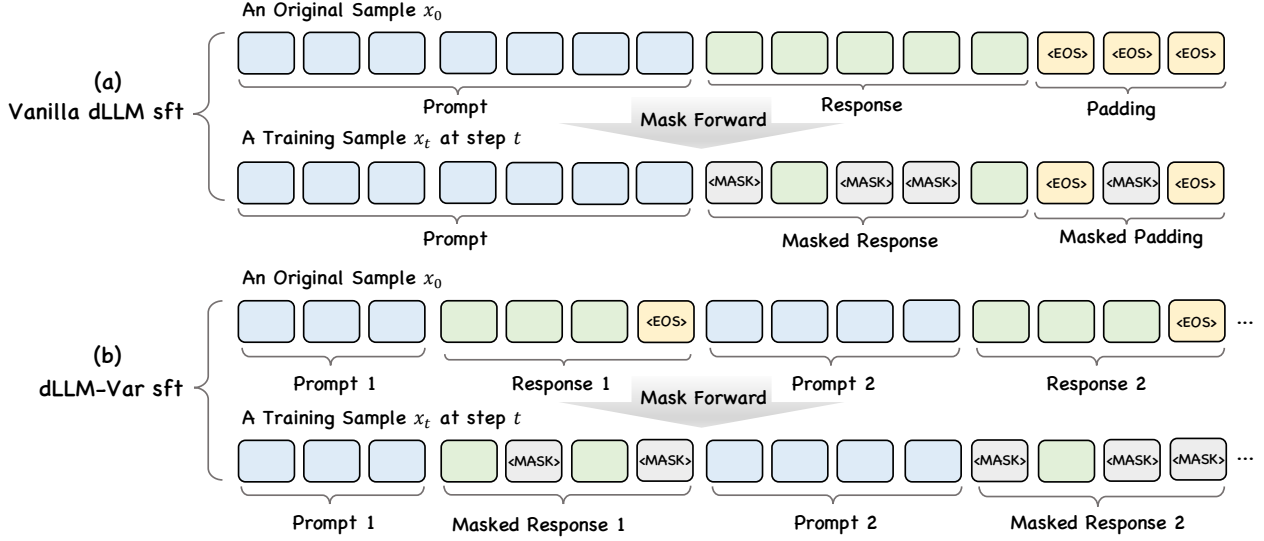


Figure 3 | During the masking forward process of dLLM-Var, tokens in the prompt are never masked. In the response section, tokens are replaced with a $\langle \text{mask} \rangle$ token based on a probability, while the final EOS token is always masked.

2.4. Block Diffusion and Attention Mechanisms

Block diffusion inference paradigms, such as Block Diffusion [13] and SDAR [8], facilitate KV cache reuse and flexible-length generation. However, their reliance on specialized attention masks imposes significant drawbacks, including a twofold increase in computational training costs and restrictive block sizes (e.g., 4-8 tokens). Our work circumvents these limitations by enabling the effects of block diffusion under a full attention mechanism, thus promoting cache efficiency without requiring specialized masks. Furthermore, this architectural choice is crucial for future advancements. Recent works on self-reflective remasking [22, 23] highlight the potential for self-correction by modifying previously generated tokens. This capability is fundamentally incompatible with block attention masks, which prevent revisiting prior blocks. In contrast, full attention provides the unrestricted access necessary for such corrective procedures. Consequently, full attention emerges not merely as an alternative, but as an essential requirement for this promising class of self-correcting dLLMs.

3. Method

In this section, we detail the core components of dLLM-Var, our proposed training framework for dLLMs. Unlike other dLLMs training methods [3, 5], dLLM-Var introduces two key innovations to enable flexible length generation and efficient KV cache reuse: (1) a fixed masking schedule for termination tokens to promote contextual awareness of sequence boundaries, and (2) multi-sample packing with full attention to enhance the model’s understanding of termination signals across unrelated contexts. These modifications enable dLLMs’ block diffusion during inference to truly understand whether the output text needs to terminate, rather than simply adding the termination tokens at the end of the first block.

3.1. Fixed Masking for Termination Tokens

During training, for a sequence $x_0 = (x_0[0], \dots, x_0[L-1])$, a noise level t is sampled uniformly from $[0, 1]$, and each token $x_0[i]$ is independently replaced with $\langle \text{MASK} \rangle$ with probability t , yielding a noisy sequence x_t . The model is then trained to predict x_t back to x_0 conditioned on the prompt.

During inference, termination tokens (e.g., EOS) must be generated by the model at appropriate positions and are not part of the input condition. To encourage this, we modify the masking schedule such that EOS tokens are always masked during training. Formally, the masking probability for position i is:

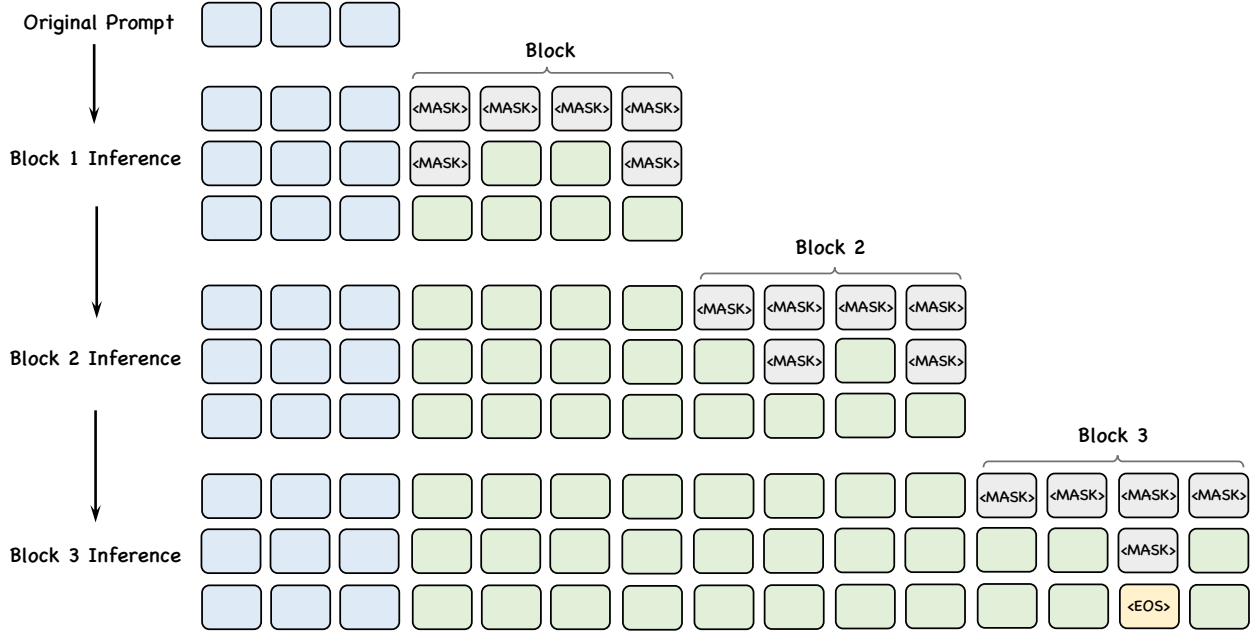


Figure 4 | The inference process of dLLM-Var. For the prompt and the already generated blocks, they will be stored in the form of a KV cache to accelerate the model’s inference.

$$p(x_t[i] = [\text{MASK}] \mid x_0[i]) = \begin{cases} t, & \text{if } x_0[i] \neq \text{EOS}, \\ 1, & \text{if } x_0[i] = \text{EOS}. \end{cases}$$

Here, x_0 denotes the original (clean) sample, x_t denotes the noisy version corresponding to the noise level t , and $i \in \{0, \dots, L-1\}$ indexes the token position within the sequence. This ensures that the model learns to output the termination token within a single conversation.

3.2. Multi-Sample Packing with Full Attention

Merely fixing the masking for EOS tokens is insufficient, as it may lead the model to naively append EOS at the end of the first block without understanding contextual semantics. To address this, we introduce multi-sample packing: unlike standard single-sample or multi-turn supervised fine-tuning (SFT) in dLLMs [1, 5], we randomly sample N dialogue pairs $\{(p^{(k)}, r^{(k)})\}_{k=1}^N$ from the dataset and concatenate them into a single training sequence of length L , separated by EOS tokens.

For the prompt parts $p^{(k)}$, which serve as generation conditions, we set the masking probability to 0. For the response parts $r^{(k)}$, masking follows the t -based schedule, and all separating EOS tokens are always masked. Crucially, we apply full attention across the entire concatenated sequence, without any attention masks.

This packing exposes the model to multiple unrelated contexts separated solely by EOS, forcing it to learn the semantic role of EOS for generations, thereby enabling it to correctly output termination token. Surprisingly, this concatenation of unrelated samples does not lead to any degradation in model performance.

3.3. Inference with Block Diffusion and KV Cache Reuse

During inference, we adopt a block-wise diffusion process akin to Block diffusion [13], but under full attention without specialized masks. Given a prompt, we progressively append blocks initialized entirely with [MASK] tokens. The next block is only added after the current block is fully unmasked and contains no EOS token, where the block size is arbitrary. Through the above method, we can achieve generation of arbitrary lengths. See Figure 4 for an illustration of the inference process.

Analogous to autoregressive models, we consider that the prompt and the already generated parts have fixed

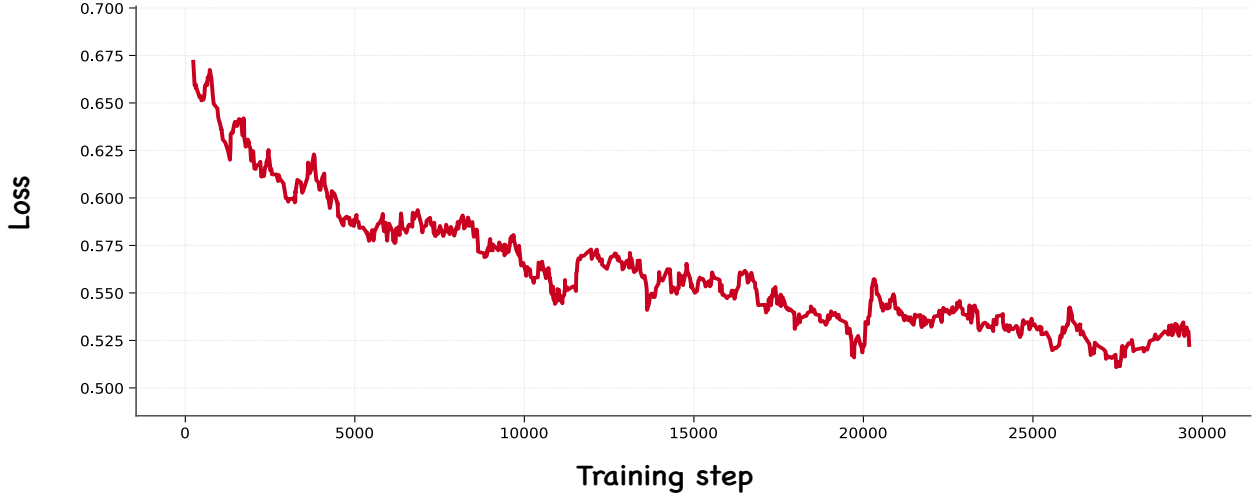


Figure 5 | Training Loss curve during our training of dLLM-Var.

semantic information under the premise that they do not need to be modified. Thus, reusing their KV caches does not compromise generation quality while substantially reducing redundant computations. Specifically, for the prompt and fully unmasked response blocks, we directly reuse their KV caches to accelerate inference. To enhance parallelism, we set a default block size of 64 tokens. Additionally, we incorporate a confidence threshold [11] of 0.9 for parallel decoding decisions, balancing between inference speed and output quality.

4. Experiments

4.1. Training Setting

We trained the dLLM-Var model using 6.4 million supervised fine-tuning (SFT) dataset pairs over 3 epochs on the DeepSpeed ZeRO-2 framework across 8 nodes with 64 GPUs. The global batch size was 1 million tokens, with each sample having a length of 8192 tokens. The learning rate was set to 1×10^{-5} , with a cosine warm-up schedule over the first 50 steps and no subsequent decay. To accelerate training, we employed FP8 mixed precision via Transformers Engines, during which no loss spikes were observed.

4.2. Evaluation Setting

During the evaluation experiments, the maximum generation length was set to 1024 for all setups, and the speedup ratio was obtained by comparing the time required to complete the model generation distribution for the evaluation tasks. All inference experiments were conducted on a single GPU, using bf16 precision, and without any additional optimizations to the PyTorch code.

4.3. Experiments Result

To validate the effectiveness of our proposed dLLM-Var, we conduct a comprehensive evaluation of its performance against several leading dLLM methodologies, including the baseline LLaDA-8B-Instruct, dLLM-Cache, and Fast dLLM. Our experiments focus on two critical metrics: inference speed and generation accuracy, with the results presented in Figure 6.

As illustrated in Figure 6 (a), dLLM-Var achieves a dramatic improvement in inference efficiency. It demonstrates a remarkable speed-up, peaking at $30.1 \times$ relative to the standard LLaDA-8B-Instruct model. This performance significantly outpaces other advanced methods such as Fast dLLM and dLLM-Cache, confirming that our approach effectively leverages parallelism and KV caching without the overhead of complex attention mechanisms.

Crucially, these substantial gains in speed do not compromise the model’s generative quality. Figure 6 (b)

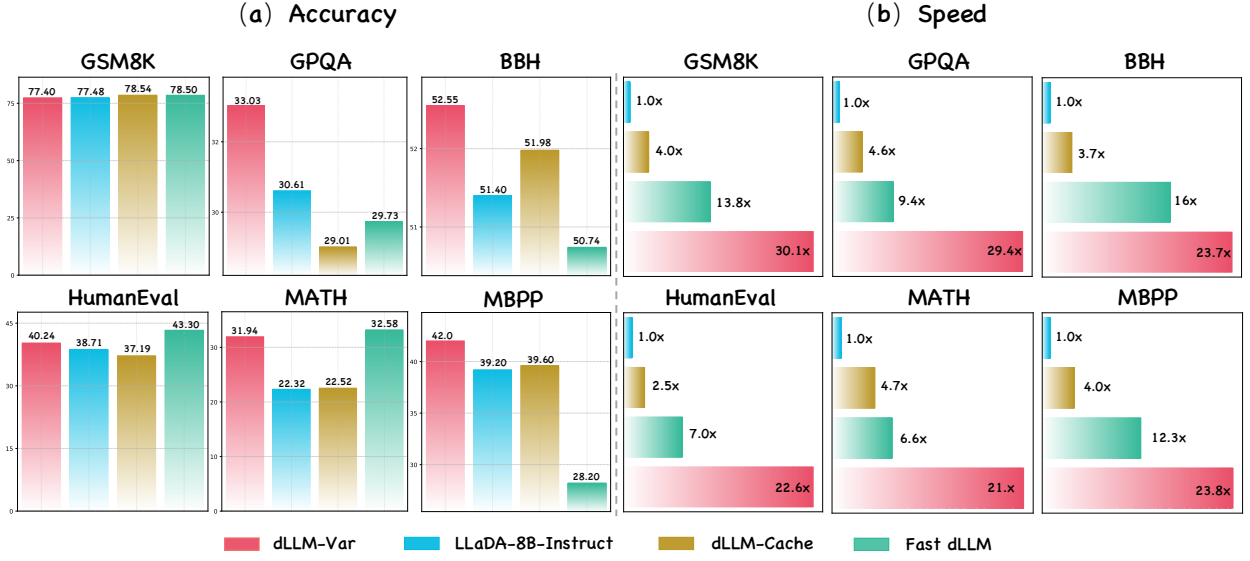


Figure 6 | **Performance comparison of dLLM-Var against baseline methods.** (a) Accuracy comparison on six standard benchmarks. The results show that dLLM-Var maintains competitive accuracy while significantly improving inference speed. (b) Speed-up ratio relative to LLaDA-8B-Instruct. dLLM-Var demonstrates substantial acceleration, achieving up to a 30.1x speed-up.

Table 1 | Accuracy and Inference Speed Across Diverse Benchmarks. We tested several models with their corresponding inference methods: Semi-ar diffusion (for LLaDA-8B-Instruct, dLLM-cache, and fast-dLLM), Pure diffusion (for LLaDA-8B-Base and dLLM-Var-pd), and our proposed Block diffusion (for dLLM-Var-bd). Speed tests were performed with a maximum generation length of 1024 using official baseline scripts. The “k-shot” designation in the benchmarks specifies the number of few-shot examples per task.

Model	Inference Method	GSM8k (5-shot)		GPQA (5-shot)		BBH (3-shot)		MATH (3-shot)		HumanEval (0-shot)		MBPP (3-shot)	
		Acc.	Speed	Acc.	Speed	Acc.	Speed	Acc.	Speed	Acc.	Speed	Acc.	Speed
Qwen3-8B	Autoregressive	89.84	17.8x	44.44	18.7x	78.40	12.3x	60.80	15.3x	65.93	9.3x	69.80	11.3x
Llama-8B-Instruct	Autoregressive	78.30	18.7x	31.90	19.2x	61.10	13.6x	29.60	18.5x	59.80	10.4x	57.60	12.3x
LLaDA-8B-Base	Pure diffusion	69.06	1.0x	31.91	1.0x	44.77	1.0x	30.84	1.0x	32.92	1.0x	40.80	1.0x
LLaDA-8B-Instruct	Semi-ar diffusion	77.48	1.0x	30.61	1.0x	51.40	1.0x	22.32	1.0x	38.71	1.0x	39.20	1.0x
LLaDA 8B Instruct	dLLM-cache	78.54	4.0x	29.01	4.6x	51.98	3.7x	22.52	4.7x	37.19	2.5x	39.60	4.0x
LLaDA-8B-Instruct	Fast dLLM	78.50	13.8x	29.73	9.4x	50.74	16x	33.20	6.6x	43.30	7.0x	28.20	12.3x
dLLM-Var-pd	Pure diffusion	77.81	1.0x	32.92	1.0x	53.71	1.0x	32.58	1.0x	39.94	1.0x	41.8	1.0x
dLLM-Var-bd (ours)	Block diffusion	77.40 (-0.41)	30.1x (+29.1x)	33.03 (+0.11)	29.4x (+28.4x)	52.55 (-1.16)	23.7x (+22.7x)	31.94 (-0.64)	21.0x (+20.0x)	40.24 (+0.30)	22.6x (+21.6x)	42.0 (+0.2)	23.8x (+22.8x)

shows the accuracy comparison across six diverse benchmarks: GSM8K, GPQA, BBH, MATH, HumanEval, and MBPP. dLLM-Var consistently delivers competitive or superior performance. For instance, it outperforms all baseline models on challenging benchmarks like BBH (52.55), GPQA (33.03), and MBPP (42.0). While its performance is on par with other methods in reasoning tasks like GSM8K, its robust results across different domains underscore that our training and inference paradigm successfully preserves, and in many cases enhances, the model’s core capabilities.

The Ability of Editing: Our experiments also highlight the self-correction capabilities of dLLM-Var. As shown in Figure 7, the model can identify and rectify errors within its own generations. When tasked with a reasoning problem, an initial version of the model’s response contained both logical and grammatical inaccuracies. Through its iterative refinement process, dLLM-Var revisited and amended this output to produce the correct and coherent final answer. This ability to self-correction demonstrates a promising capacity for reasoning about and improving its own output, marking a crucial step towards more reliable and accurate generative models.

Table 2 | Accuracy and speed-up comparison for dLLM-Var with different methods and lengths on the GSM8K and MBPP tasks. For Pure diffusion, the length represents the fixed generation quota.

Inference Method	Generation Length	GSM8K (5-shot)		MBPP (3-shot)	
		Accuracy	Speed	Accuracy	Speed
Pure diffusion	64	12.35	21.4×	38.40	19.8×
	128	51.93	10.1×	42.00	9.1×
	256	74.98	4.8×	42.00	4.4×
	512	75.96	2.3×	43.20	2.1×
	1024	77.81	1.0×	41.80	1.0×
Block diffusion	-	77.40	30.1×	42.00	23.8×

Prompt: Janet’s ducks lay 16 eggs per day. She eats three for breakfast every morning and bakes muffins for her friends every day with four. She sells the remainder at the farmers' market daily for \$2 per fresh duck egg. How much in dollars does she make every day at the farmers' market?

Before editing

Janet has 16 eggs per day. She eats 3 eggs for breakfast ,
so she has 16 - 3 = 13 eggs left . She uses 4 eggs
for her muffins , so she has 13 - 4 = 9 eggs left . She
sells the remaining 8 eggs at the farmers' market for \$2 per
egg , so he makes 8 * \$2 = \$16 per day . \n ### 1 6

After editing

Janet has 16 eggs per day. She eats 3 eggs for breakfast ,
so she has 16 - 3 = 13 eggs left . She uses 4 eggs
for her muffins , so she has 13 - 4 = 9 eggs left . She
sells the remaining 9 eggs at the farmers' market for \$2 per
egg , so she makes 9 * \$2 = \$18 per day . \n ### 1 8

Figure 7 | Demonstration of dLLM-Var’s Self-Correction Capability. It illustrates the model’s ability to refine its initial output. In the "Before editing" phase, the model makes both calculation error and grammatical error. In the "After editing" phase, the model corrects these mistakes, adjusting the calculation to the correct 9 eggs and changing the pronoun to "she", resulting in the correct final answer of 18.

5. Conclusion and Future Works

dLLM-Var revolutionizes diffusion-based LLMs by shattering fixed-length barriers and unlocking seamless KV cache reuse through an elegant training method: fixed EOS masking and multi-sample packing. With this work, open-source dLLMs transcend academic novelty, unlocking real-world industrial viability. While dLLM-Var marks a significant step towards making dLLMs practical, several exciting avenues for future research remain. Our current approach treats previously generated blocks as fixed, but this unidirectional process forgoes the opportunity for self-correction based on newly generated context. A promising direction is to explore mechanisms for iterative refinement and editing of generated content. As the model gains a richer contextual understanding from subsequent blocks, it could re-evaluate and selectively re-generate portions of earlier output to correct inaccuracies and improve coherence. The full attention mechanism employed in dLLM-Var is a key enabler for this research, allowing unrestricted access to all tokens. Future work could focus on developing efficient algorithms to decide "when" and "what" to edit, balancing the trade-off between quality gains and computational cost. Such advancements could unlock a new level of performance and reliability for dLLMs, further closing the gap with their autoregressive counterparts.

References

- [1] 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.
- [2] 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. *arXiv preprint arXiv:2505.16933*, 2025.
- [3] 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. *arXiv preprint arXiv:2505.19223*, 2025.
- [4] Yuxuan Song, Zheng Zhang, Cheng Luo, Pengyang Gao, Fan Xia, Hao Luo, Zheng Li, Yuehang Yang, Hongli Yu, Xingwei Qu, et al. Seed diffusion: A large-scale diffusion language model with high-speed inference. *arXiv preprint arXiv:2508.02193*, 2025.
- [5] Jiacheng Ye, Zhihui Xie, Lin Zheng, Jiahui Gao, Zirui Wu, Xin Jiang, Zhenguo Li, and Lingpeng Kong. Dream 7b, 2025.
- [6] Zhihui Xie, Jiacheng Ye, Lin Zheng, Jiahui Gao, Jingwei Dong, Zirui Wu, Xueliang Zhao, Shansan Gong, Xin Jiang, Zhenguo Li, et al. Dream-coder 7b: An open diffusion language model for code. *arXiv preprint arXiv:2509.01142*, 2025.
- [7] 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. *arXiv preprint arXiv:2506.20639*, 2025.
- [8] Shuang Cheng, Yihan Bian, Dawei Liu, Yuhua Jiang, Yihao Liu, Linfeng Zhang, Wenghai Wang, Qipeng Guo, Kai Chen, Biqing Qi*, and Bowen Zhou. Sdar: A synergistic diffusion–autoregression paradigm for scalable sequence generation, 2025.
- [9] Chengyue Wu, Hao Zhang, Shuchen Xue, Shizhe Diao, Yonggan Fu, Zhijian Liu, Pavlo Molchanov, Ping Luo, Song Han, and Enze Xie. Fast-dllm v2: Efficient block-diffusion llm, 2025.
- [10] Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière, David Lopez-Paz, and Gabriel Synnaeve. Better & faster large language models via multi-token prediction. *arXiv preprint arXiv:2404.19737*, 2024.
- [11] 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.
- [12] Xu Wang, Chenkai Xu, Yijie Jin, Jiachun Jin, Hao Zhang, and Zhijie Deng. Diffusion llms can do faster-than-ar inference via discrete diffusion forcing. *arXiv preprint arXiv:2508.09192*, 2025.
- [13] Marianne Arriola, Aaron Gokaslan, Justin T Chiu, Zhihan Yang, Zhixuan Qi, Jiaqi Han, Subham Sekhar Sahoo, and Volodymyr Kuleshov. Block diffusion: Interpolating between autoregressive and diffusion language models. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [14] Zigeng Chen, Gongfan Fang, Xinyin Ma, Ruonan Yu, and Xinchao Wang. dparallel: Learnable parallel decoding for dllms. *arXiv preprint arXiv:2509.26488*, 2025.
- [15] 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. *arXiv preprint arXiv:2506.06295*, 2025.
- [16] Yuchu Jiang, Yue Cai, Xiangzhong Luo, Jiale Fu, Jiarui Wang, Chonghan Liu, and Xu Yang. d2cache: Accelerating diffusion-based llms via dual adaptive caching. *arXiv preprint arXiv:2509.23094*, 2025.
- [17] Qingyan Wei, Yaojie Zhang, Zhiyuan Liu, Dongrui Liu, and Linfeng Zhang. Accelerating diffusion large language models with slowfast: The three golden principles. *arXiv preprint arXiv:2506.10848*, 2025.

- [18] Yuxin Ma, Lun Du, Lanning Wei, Kun Chen, Qian Xu, Kangyu Wang, Guofeng Feng, Guoshan Lu, Lin Liu, Xiaojing Qi, Xinyuan Zhang, Zhen Tao, Haibo Feng, Ziyun Jiang, Ying Xu, Zenan Huang, Yihong Zhuang, Haokai Xu, Jiaqi Hu, Zhenzhong Lan, Junbo Zhao, Jianguo Li, and Da Zheng. dinfer: An efficient inference framework for diffusion language models. *arXiv preprint arXiv:2510.08666*, 2025.
- [19] Zirui Wu, Lin Zheng, Zhihui Xie, Jiacheng Ye, Jiahui Gao, Yansong Feng, Zhenguo Li, Victoria W., Guorui Zhou, and Lingpeng Kong. Dreamon: Diffusion language models for code infilling beyond fixed-size canvas, 2025.
- [20] Jinsong Li, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Jiaqi Wang, and Dahua Lin. Beyond fixed: Variable-length denoising for diffusion large language models. *arXiv e-prints*, pages arXiv–2508, 2025.
- [21] Jaeyeon Kim, Lee Cheuk-Kit, Carles Domingo-Enrich, Yilun Du, Sham Kakade, Timothy Ngatiaoco, Sitan Chen, and Michael Albergo. Any-order flexible length masked diffusion. *arXiv preprint arXiv:2509.01025*, 2025.
- [22] Zemin Huang, Yuhang Wang, Zhiyang Chen, and Guo-Jun Qi. Don’t settle too early: Self-reflective remasking for diffusion language models. *arXiv preprint arXiv:2509.23653*, 2025.
- [23] Guanghan Wang, Yair Schiff, Subham Sahoo, and Volodymyr Kuleshov. Remasking discrete diffusion models with inference-time scaling. *arXiv preprint arXiv:2503.00307*, 2025.