

# Efficient-DLM: From Autoregressive to Diffusion Language Models, and Beyond in Speed

Yonggan Fu\*, Lexington Whalen\*<sup>1</sup>, Zhifan Ye<sup>1</sup>, Xin Dong, Shizhe Diao, Jingyu Liu<sup>2</sup>, Chengyue Wu<sup>3</sup>, Hao Zhang, Enze Xie, Song Han<sup>4</sup>, Maksim Khadkevich, Jan Kautz, Yingyan (Celine) Lin<sup>1</sup>, Pavlo Molchanov

Diffusion language models (dLMs) have emerged as a promising paradigm enabling parallel generation, but their learning efficiency lags behind that of autoregressive (AR) language models when trained from scratch. To this end, we study AR-to-dLM conversion, which transforms pretrained AR models into efficient dLMs that excel in speed while preserving AR models' task accuracy. We achieve this by identifying limitations in the attention patterns and objectives of existing AR-to-dLM methods and then proposing methodologies and actionable insights for scalable AR-to-dLM conversion. Specifically, we first systematically compare different attention patterns and find that maintaining pretrained AR weight distributions is key to effective AR-to-dLM conversion. Accordingly, we introduce a continuous pretraining scheme with a block-wise attention pattern, which remains causal across blocks with bidirectional modeling within each block. We find that, in addition to block-wise attention's known benefit of enabling KV caching, its block-wise causality better preserves pretrained AR models' weight distributions, leading to a win-win in accuracy and efficiency. Second, to mitigate the training-test gap in mask token distributions (uniform vs. highly left-to-right), we propose a position-dependent token masking strategy that assigns higher masking probabilities to later tokens during training to better mimic test-time behavior. Leveraging this framework, we conduct extensive studies of dLMs' attention patterns, training dynamics, and other design choices. These studies lead to the Efficient-DLM model family, which outperforms state-of-the-art AR models and dLMs in accuracy-throughput trade-offs; for example, our Efficient-DLM 8B achieves +5.4%/+2.7% higher accuracy with  $4.5\times/2.7\times$  higher throughput compared to Dream 7B and Qwen3 4B, respectively.

**Models on Hugging Face:** [Efficient-DLM-4B](#) | [Efficient-DLM-8B](#)

## 1. Introduction

The success of large language models (LLMs) has been largely driven by autoregressive (AR) modeling, where tokens are generated sequentially left to right. Despite strong benchmark performance, AR models are constrained by token-by-token decoding, which limits generation throughput, especially in memory-bounded scenarios (e.g., small batch sizes) where hardware utilization is low.

To overcome the sequential bottleneck of AR decoding, diffusion language models (dLMs) [1, 2, 3, 4] have recently emerged as an alternative paradigm. By leveraging iterative denoising steps, dLMs enable parallel, non-autoregressive generation and hold promise for higher throughput. However, despite their conceptual appeal, most existing dLMs have not delivered faster speed than AR models in practice [3, 4], due to the limited compatibility with key-value (KV) caching and the limited parallelism during decoding. Although pioneering works [5, 6, 7] demonstrate potential speed-up on small-scale models (e.g., 110M [5]) with limited downstream accuracy, successful scaling of dLMs to larger model sizes has been restricted by

prohibitive training costs [8]. This is because AR models learn only left-to-right modeling, while dLMs learn all possible permutations [9], which is more difficult and requires longer training.

This work leverages pretrained AR models for initialization and systematically explores how to continuously pretrain them into dLMs that achieve high generation speed while preserving task accuracy. The key insight is that, with an appropriate training scheme in terms of attention patterns and objectives, pretrained AR models can be converted into faster dLMs that support parallel decoding with KV cache at low training cost (on the order of 10B tokens), and extended continuous training (on the order of 100B tokens) enables more aggressive parallel generation.

We achieve this by identifying limitations in the attention patterns and objectives of existing AR-to-dLM methods and proposing a continuous pretraining scheme that features a block-wise attention pattern [5] and position-dependent token masking. Specifically, our extensive study of attention patterns shows that a block-wise attention pattern, beyond its known benefit of enabling KV caching, better preserves the weight

\* Co-first author; additional affiliation: <sup>1</sup>Georgia Tech, <sup>2</sup>University of Chicago, <sup>3</sup>University of Hong Kong, <sup>4</sup>MIT.  
© 2026 NVIDIA. All rights reserved.

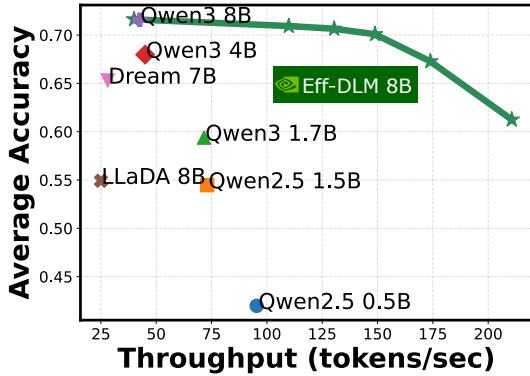


Figure 1 | Benchmarking the accuracy–throughput trade-offs between Efficient-DLM 8B and SOTA AR/dLMs, with accuracy averaged over 12 tasks across math, coding, and commonsense reasoning.

distributions of pretrained AR models than the fully bidirectional training adopted in prior work [3, 4]. This preservation is crucial for effective AR-to-dLM conversion while maintaining task accuracy. These findings provide a practical guideline for the community: block-wise attention (with each block conditioned on clean context during training) is a preferred choice for AR-to-dLM conversion, achieving a win–win in accuracy and efficiency.

In addition, we identify a mismatch between training-time uniform token masking and test-time confidence-based token sampling. To bridge this gap and improve downstream accuracy, we propose a position-dependent token masking strategy. This approach builds on the observation that dLMs retain a left-to-right generation tendency due to the autoregressive nature of language, and thus incorporates the prior that, as the input becomes less corrupted (i.e., closer to complete denoising), more tokens should be masked toward the end of each block.

Furthermore, under this training framework, we analyze attention patterns, training dynamics, and other design choices for scalable AR-to-dLM conversion, and introduce the Efficient-DLM model family, which outperforms both AR and dLM baselines with improved accuracy–throughput trade-offs. For example, as shown in Fig. 1, our Efficient-DLM 8B maintains comparable (slightly better) accuracy than Qwen3 8B and achieves +5.4%/+2.7% higher accuracy with  $4.5\times/2.7\times$  higher throughput compared to Dream 7B and Qwen3 4B, respectively.

We expect these findings to provide practical guidelines for realizing dLMs’ promise of faster, more efficient generation, and to inspire new dLM paradigms. The key takeaways and insights from this work are summarized here.

### Takeaways for converting pretrained AR models into faster dLMs

- *Attention pattern is key to AR-to-dLM conversion:* Continuous training with block attention better preserves pretrained AR models’ abilities than fully bidirectional modeling; Conditioning each corrupted block on clean context is essential for model accuracy.
- *Training block size matters:* Too-small block sizes lack sufficient context for denoising, while too-large block sizes induce excessive corruption and weight changes.
- *Evaluation block size matters:* Training with proper block sizes can generalize well to other evaluation block sizes; Larger evaluation block sizes generally provide more opportunities for parallel decoding.
- *Whether to preserve the token shift of AR models:* We find it unnecessary and potentially harmful.
- *Left-to-right generation tendency:* dLLMs still exhibit this tendency during parallel generation, and mimicking it in training can boost generation quality.
- *Training dynamics:* Likelihood estimation improves steadily with training, allowing for more aggressive parallel decoding.

## 2. Efficient-DLM: Attention Pattern Analysis

### 2.1. Analyzing Different Attention Patterns

**Fully bidirectional attention in existing works.** Existing works that transform AR models into dLMs [10, 4] adopt fully bidirectional modeling, i.e., the entire sequence is randomly corrupted and all tokens are visible to each other, as shown in Fig. 2 (a) and (b). This training scheme suffers from the following drawbacks: (1) fully bidirectional attention increases the difficulty of applying KV caching; (2) the context is overly corrupted, particularly for later tokens, which increases training difficulty; (3) the fully bidirectional attention pattern diverges from the causality of the AR initialization, resulting in larger weight drifts from pretrained AR models.

**Block-wise attention.** In light of these limitations, an enhanced attention pattern is the block-wise attention shown in Fig. 2 (c), which remains causal across blocks and adopts bidirectional modeling within each block. This design enables the use of KV caching at test time, and its block-wise causality is closer to the token-wise causality of pretrained AR models, potentially better preserving their abilities.

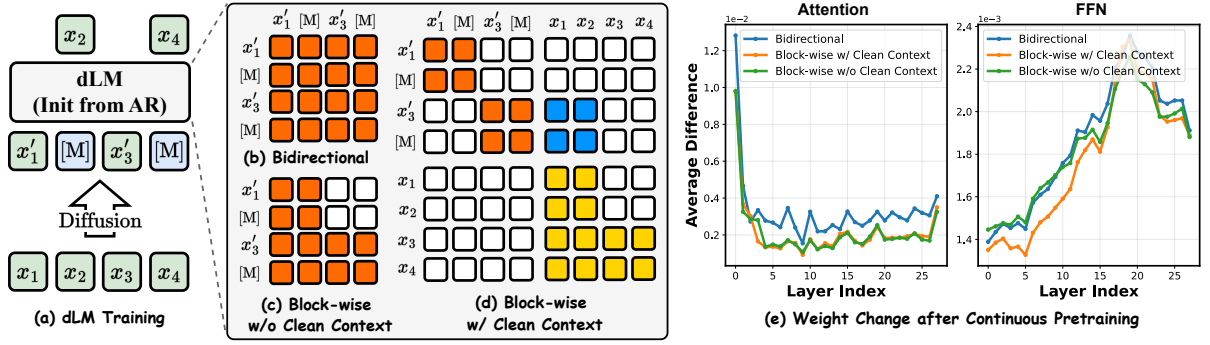


Figure 2 | Visualizing continuous pretraining of dLMs with different attention patterns from pretrained AR models. (b) and (c) show bidirectional attention and block-wise attention without clean context, respectively, using a block size of 2 as an example. (d) illustrates the block-wise attention with clean context (using a block size of 2 as an example), where  $\blacksquare$  denotes attention among noisy tokens,  $\blacksquare$  denotes attention from noisy tokens to clean-context tokens, and  $\blacksquare$  denotes attention within the clean context. (e) shows weight changes in the attention and feed-forward network (FFN) layers after training under these three attention patterns.

**Block-wise attention with each block conditioned on clean context.** The drawback of the block-wise attention in Fig. 2 (c) is that it may cause a training–test gap: when performing block-wise decoding at test time, the context preceding a noisy block has already been decoded without mask tokens; however, the attention pattern in Fig. 2 (c) cannot ensure this during training, as the context of each block still contains mask tokens. To resolve this training–test gap, our work employs block-wise attention with each block conditioned on clean context, as shown in Fig. 2 (d), following [5]. This ensures that each block is conditioned only on clean context during training, mimicking the block-wise decoding process at test time, where all previous blocks are fully completed without mask tokens. This is achieved by concatenating the noisy tokens and the clean tokens as dLM inputs and applying the special attention mask shown in Fig. 2 (d). Such an attention pattern allows for seamless use of the KV cache for improved efficiency, constrains corruption within each block to maintain a cleaner context, and preserves block-wise causality, thereby mitigating weight drift and better inheriting the capabilities of the original AR model.

More formally, let  $\mathbf{x} = (x_1, \dots, x_L)$  be a sequence partitioned into  $B$  contiguous blocks  $\mathbf{x}^b$  of length  $L' = L/B$ , and let  $q(\tilde{\mathbf{x}}_t^b | \mathbf{x}^b)$  denote the corruption process that produces the noisy input  $\tilde{\mathbf{x}}_t^b$  at noise level  $t \in (0, 1]$ . Our training objective is as follows:

$$\mathcal{L}(\theta) = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{\tilde{\mathbf{x}}_t^b \sim q(\cdot | \mathbf{x}^b)} \left[ -\frac{1}{t} \sum_{b=1}^B \log p_\theta(\mathbf{x}^b | \tilde{\mathbf{x}}_t^b, \mathbf{x}^{<b}) \right], \quad (1)$$

where  $p_\theta(\mathbf{x}^b | \tilde{\mathbf{x}}_t^b, \mathbf{x}^{<b})$  denotes the denoising of the  $b$ -th block based on the corrupted input  $\tilde{\mathbf{x}}_t^b$  and the clean context  $\mathbf{x}^{<b}$ . As a side note, the block-wise attention variant in Fig. 2 (c) without using clean context can be formulated by replacing the clean

context  $\mathbf{x}^{<b}$  in Eq. 1 with the corrupted context  $\tilde{\mathbf{x}}^{<b}$ .

Unlike prior block diffusion methods trained from scratch [5], we initialize  $\theta$  from a pretrained AR model that is trained with the autoregressive loss  $\mathcal{L}_{\text{AR}}(\theta) = -\sum_{\ell=1}^L \log p_\theta(x_\ell | x_{<\ell})$ , and then adapt it through continuous pretraining using the loss in Eq. 1. This initialization allows for rapid AR-to-dLM conversion, which requires (1) adapting weights to new attention patterns and (2) avoiding large weight drifts to better preserve the original model’s ability.

## 2.2. Comparison of Attention Patterns

We study the impact of three key design factors on attention patterns in AR-to-dLM conversion: (1) fully bidirectional vs. block-wise; (2) whether to keep clean context: using clean context  $\mathbf{x}^{<b}$  or corrupted context  $\tilde{\mathbf{x}}^{<b}$  in Eq. 1; and (3) whether to perform token shift, i.e., predicting the next token as in AR models or directly predicting the mask tokens themselves. Previous works [10, 4] find that preserving token shift when initializing from AR models is beneficial, and we revisit this design choice under more advanced training schemes.

**Settings.** We adopt Qwen2.5 1.5B [11] as the AR initialization and perform continuous pretraining for 50B tokens on a mixed dataset comprising [12, 13, 14]. For block-wise training, we adopt a block size of 16, and provide further analysis on block sizes in Sec. 2.3. The initial learning rate is set to  $1e-5$  and decayed to  $3e-6$  using a cosine schedule with the AdamW optimizer. An analysis of the learning rate is provided in Appendix D. We evaluate downstream task accuracy on six generation tasks, including HumanEval, HumanEval Plus, MBPP, MBPP Plus, GSM8K, and Minerva Math, using lm-evaluation-harness [15].

**Importance of block-wise attention.** As shown

Table 1 | Comparing different dLM training schemes on Qwen2.5 1.5B. Row (a) shows the accuracy of the original Qwen2.5 1.5B. Row (b) presents the training scheme of Dream [4]. Row (g) shows the identified best scheme with block-wise attention, clean context, and no token shift.

Row ID	Attn Pattern	Clean Context	Token Shift	KV Cache	Human -Eval	Human -Eval Plus	MBPP	MBPP Plus	GSM8K	Minerva Math	Avg
a	AR	-	✓	✓	36.59	29.88	43.6	59.52	54.74	26.40	41.79
b	Bidirectional	-	✓	✗	15.85	12.20	16.2	24.34	28.96	11.08	18.10
c	Bidirectional	-	✗	✗	19.51	15.24	17.2	24.34	28.20	11.22	19.29
d	Block-wise	✗	✓	✓	31.10	25.61	23.6	36.77	38.44	13.88	28.23
e	Block-wise (2×)	✗	✓	✓	26.22	22.56	26.0	42.33	36.69	12.56	27.73
f	Block-wise	✓	✓	✓	38.41	33.54	33.0	48.68	51.48	21.04	37.69
g	Block-wise	✓	✗	✓	39.02	34.76	34.0	48.15	52.99	21.56	38.41

in Tab. 1, compared to bidirectional attention in Row (c), block-wise attention (even without clean context) in Row (d) can boost average accuracy by 8.94%. When combined with other best practices, i.e., conditioning on clean context and removing token shift in Row (g), block-wise attention improves the average accuracy over bidirectional attention by 19.12%.

This implies that block-wise attention better preserves block-wise causality and thus maintains the pretrained AR model’s abilities more effectively than bidirectional attention, in addition to the benefit of native KV caching. Furthermore, visualizations of weight changes after continuous pretraining in Fig. 2 (e) show that bidirectional attention leads to larger weight drifts from pretrained weights in both the attention and FFN layers, ultimately causing larger accuracy drops.

**The impact of clean context.** Based on the comparison between Rows (d) and (f) in Tab. 1, conditioning each block on clean context during training is critical, yielding a 9.46% accuracy improvement over using noisy context, where both adopt block-wise attention. We further train the noisy-context case with a doubled token budget, i.e., extend Row (d) to Row (e), to account for the increased sequence length caused by concatenating noisy and clean tokens in the setting of Row (f). However, comparing Rows (e) and (f) in Tab. 1 shows that doubling training tokens on corrupted context cannot effectively recover the accuracy, whereas training on fewer tokens with clean context yields substantially higher accuracy.

In addition, Fig. 2 (e) shows that training with block-wise attention without clean context leads to larger weight drifts in FFN layers compared to training with clean context.

**Whether to perform token shift.** We find that token shift is unnecessary, and its removal consistently improves accuracy across settings, as evidenced by the comparison between Rows (b) & (c) and Rows (f) & (g) in Tab. 1. This indicates that (1) the token shift inherent to AR models can be easily adapted

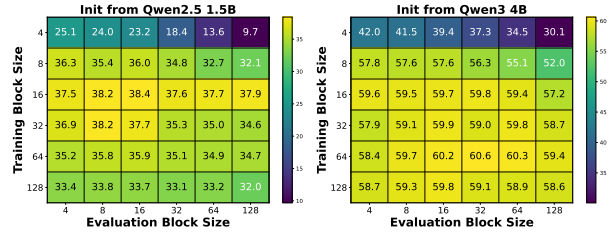


Figure 3 | The average accuracy achieved by different training–evaluation block size pairs.

into the no-token-shift setting, and (2) predicting the mask token itself (without token shift) is easier than predicting the next token of a masked position. We hypothesize that the latter is harder because the model must handle two tasks simultaneously: inferring the mask token and predicting the following token.

**Takeaways:** When continuously pretraining from an AR model, a block-wise attention pattern with clean context and without token shift emerges as a promising training scheme to deliver dLMs.

Based on this takeaway, we adopt this scheme by default in the following study.

### 2.3. Analysis of the Optimal Block Sizes

Building on the best training scheme in Sec. 2.2, the next question is the optimal block size for training and evaluation. Intuitively, larger context sizes provide richer context with more visible future tokens, but at the same time introduce more corruption, i.e., the last tokens in a block encounter noisier past context. This makes it critical to select a proper block size that balances both aspects. We study the impact of training and evaluation block sizes in this subsection.

**Settings.** We perform continuous pretraining on top of Qwen2.5 1.5B [11] and Qwen3 4B [16] for 50B and 25B tokens, respectively, using different training block sizes [4, 8, 16, 32, 64, 128]. Other configura-

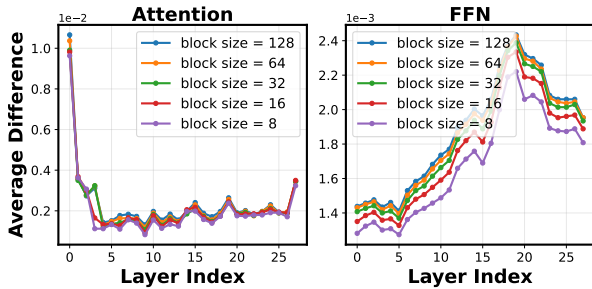


Figure 4 | The weight changes in attention and FFN layers after training with different block sizes.

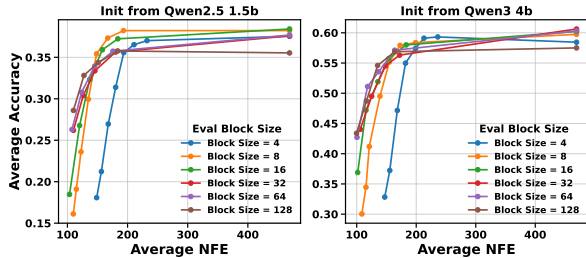


Figure 5 | The average NFE accuracy on six tasks of different evaluation block sizes under varying NFEs.

tions are the same as in Sec. 2.2. We then evaluate each trained model with different evaluation block sizes and report the average accuracy across six generation tasks (HumanEval, HumanEval Plus, MBPP, MBPP Plus, GSM8K, and Minerva Math) for each training–evaluation block size pair in Fig. 3.

**Observations on training block sizes.** As shown in Fig. 3, we observe that (1) For both model scales, too-small training block sizes generally lead to suboptimal accuracy, potentially because the context is not sufficiently rich to predict the corruptions. (2) Larger-scale models are more tolerant of larger training block sizes, which introduce more corruptions but also provide richer context. In contrast, small-scale models have a more notable sweet-spot training block size (e.g., 16 for Qwen2.5 1.5B), beyond which the more corrupted context leads to degraded accuracy. (3) Training with an appropriate block size can transfer well to other evaluation block sizes. This can be inferred from the attention map in Fig. 2 (d), where the model trained with a single block size can see varying numbers of tokens participating in the attention mechanism. This differs from the fully bidirectional attention in Fig. 2 (b), where the model always sees the same number of tokens participating in attention, necessitating additional techniques such as random sequence length truncation [3] to generalize to other sequence lengths.

We also visualize the weight changes before and after continuous pretraining with different block sizes on top of Qwen2.5 1.5B in Fig. 4. In general, larger training block sizes lead to larger weight changes, and

there exists a sweet-spot block size that yields the best downstream task accuracy. This indicates a trade-off between maintaining original abilities and adapting to new attention patterns, and the sweet-spot block size balances both aspects.

**Takeaways:** There exists a sweet-spot training block size: too-small block sizes lack sufficient context, while too-large block sizes induce excessive corruption and weight changes.

**Observations on evaluation block sizes.** To understand the impact of evaluation block sizes for a trained model, we adopt confidence-based sampling [17] with different confidence thresholds to control the number of function evaluations (NFEs) [18]. The lower the NFE, the more tokens are generated in parallel. We adopt diffusion Qwen2.5 1.5B and Qwen3 4B, trained with the best block sizes from Fig. 3 (16 and 64, respectively), and evaluate them with different block sizes. As shown in Fig. 5, we observe that (1) larger evaluation block sizes generally lead to higher accuracy when performing more aggressive token generation with lower NFEs. We assume this is because, under the same number of denoised tokens per step, larger block sizes provide more flexible positions and greater opportunities for parallel token generation; and (2) with larger NFEs, there is no clearly optimal evaluation block size, and moderate block sizes generally yield comparable results.

**Takeaways:** Although a proper training block size can generalize to other evaluation block sizes, larger evaluation block sizes favor more aggressive parallel token generation.

### 3. Efficient-DLM: Position-dependent Token Masking

#### 3.1. The Training-Test Gap in Token Masking

Existing dLMs [3, 4, 5] typically adopt uniform token masking, where mask tokens are randomly sampled from a uniform distribution based only on the noise level  $t$ , independent of token positions. However, we find that at inference time, when performing confidence-based sampling [3, 4], the denoised tokens are not uniformly distributed; instead, they show a clear left-to-right tendency.

To demonstrate this, we visualize the average number of denoising steps required at each token position in a block on the GSM8K dataset using the trained diffusion Qwen2.5 1.5B model from Sec. 2.2 in Fig. 6

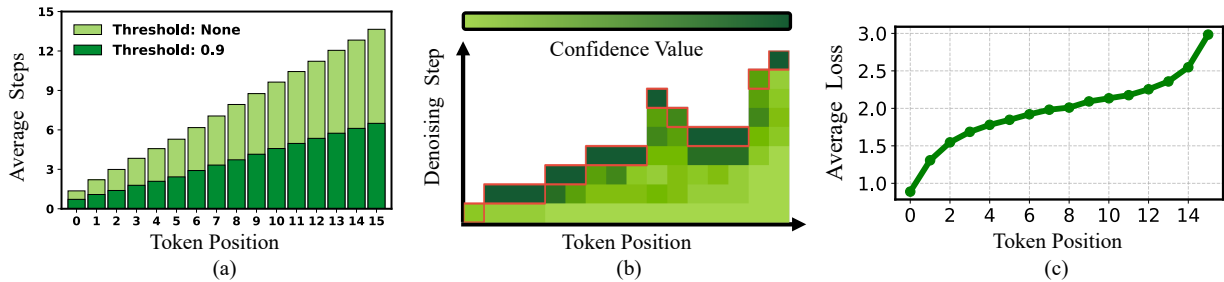


Figure 6 | (a) The average number of denoising steps required at each token position on GSM8K using diffusion Qwen2.5 1.5B with two different confidence thresholds, where "None" denotes one token per step. (b) The confidence distribution within a block across different denoising steps for an example from GSM8K, with red boxes marking tokens that are decoded. (c) The average loss at each token position within a block of size 16.

(a). Specifically, we visualize two cases, with and without parallel token generation, using a confidence threshold [17], where tokens with confidence surpassing this threshold are decoded. We observe that the average number of denoising steps increases with the positions in a block, exhibiting a notable left-to-right tendency due to the autoregressive nature of language.

As a concrete example, we also show the confidence distribution within a block across different denoising steps for one example from GSM8K in Fig. 6 (b), with red boxes marking tokens that are decoded and finalized. We can see that tokens tend to have higher confidence scores once their neighboring tokens have been decoded and are generally decoded from left to right. In other words, as the denoising process approaches completion, mask tokens are more likely to appear near the block end. As such, uniform token masking during training and confidence-based sampling during inference create a training–test gap.

In addition, we visualize the average loss at each token position within a block in Fig. 6 (c), averaged over all blocks for 200 samples. We observe that the later mask tokens in a block are generally harder cases with larger losses due to more corrupted context, potentially requiring more learning. This also indicates that a more strategic token masking scheme that also considers token position is desirable.

### 3.2. Our Token Masking Strategy

To demonstrate the importance of mitigating the identified training–test gap, we propose the concept of position-dependent token masking. Specifically, for a sequence  $\mathbf{x} = (x_1, \dots, x_L)$  and a given noise level  $t$ , conditioned on  $t$  and the relative token position  $i \in [L']$  within one block, the masking probability of each token position is set as

$$w_i(t) = \exp[\beta(1-t)i], \quad (2)$$

where  $\beta \geq 0$  is a hyperparameter controlling the strength of the positional bias. Specifically,  $\beta = 0$

leads to uniform sampling, and larger  $\beta$  indicates a stronger positional bias. The set of mask tokens is drawn from this distribution by normalizing the weights and then performing Gumbel-top- $k$  sampling [19], where  $k = \lfloor tL' \rfloor$  is the per-block mask token count.

When  $t \rightarrow 0$ , corresponding to the end of denoising,  $w_i(t)$  assigns larger weights to later tokens, i.e., mask tokens are more likely to appear near the block end, echoing the test-time pattern in Sec. 3.1. When  $t \rightarrow 1$ , corresponding to noisier inputs,  $w_i(t)$  becomes more uniform, i.e., mask tokens are sampled more uniformly. Such a position-dependent token masking narrows the training–test gap and increases the masking probability of later tokens, i.e., the harder cases with larger losses. We also note that related work [20] leverages the AR nature of language and designs a left-to-right noise schedule for continuous dLMs, whereas our work focuses on discrete dLMs.

### 3.3. Comparison of Token Masking Schemes

**Settings.** We apply position-dependent token masking with different  $\beta$  to the training of diffusion Qwen3 4B with block size 64 on 25B tokens. In practice, instead of directly setting  $\beta$ , we parameterize the positional prior using a half-life ratio  $\lambda = \ln 2 / (\beta L') \in (0, 1]$ . The half-life ratio  $\lambda$  is the fraction of a block length over which, under maximal tilt  $t \rightarrow 0$ , the positional weight changes by a factor of two. Thus, the lower the value of  $\lambda$ , the stronger the positional prior.

We compare position-dependent token masking with different  $\lambda$  values against uniform token masking (i.e.,  $\lambda \rightarrow \infty$ ) and right-to-left masking (i.e.,  $\lambda \rightarrow 0$ ), which always masks the rightmost  $k$  tokens. The average masking probability of each position within a block throughout training is shown in Fig. 7. The average accuracy on the six generation tasks in Sec. 2.3, under different parallel decoding settings in terms of tokens per forward (TPF), i.e., the number of decoded

Setting	TPF=1	TPF=2.8	TPF=4	TPF=5.6	Avg Diff
uniform	60.27	57.12	51.11	33.99	-
right-to-left	38.21	27.55	18.60	14.13	-
$\lambda=0.25$	60.41 (+0.14)	58.56 (+1.44)	50.97 (-0.14)	34.55 (+0.56)	+0.50
$\lambda=0.1$	62.02 (+1.75)	58.79 (+1.67)	53.75 (+2.64)	38.37 (+4.38)	+2.61
$\lambda=0.05$	60.51 (+0.24)	57.68 (+0.56)	53.06 (+1.95)	37.38 (+3.39)	+1.54

Table 2 | Comparing token masking schemes based on the average accuracy across six generation tasks under varying parallel decoding settings, measured in tokens per forward (TPF).

tokens per denoising step using confidence-based sampling [17], is presented in Tab. 2.

**Observations.** As shown in Tab. 2, we observe that (1) progressively increasing positional priors with lower  $\lambda$  leads to improved average accuracy; (2) positional priors are particularly beneficial under more aggressive parallel decoding settings, with up to a 4.38% average accuracy improvement; and (3) positional priors should not be blindly increased, as the extreme case of fully right-to-left masking leads to poor results, likely because the model is forced to train only on the hard cases at the block end without learning to exploit the bidirectional context.

These experiments indicate that positional priors are helpful but must be introduced properly. The key contribution of our method here is to introduce this design factor, and we hope it can inspire more advanced and automated schemes in the future.

**Takeaways:** dLMs exhibit a left-to-right tendency during parallel generation due to the autoregressive nature of language, and mimicking this tendency in training can boost generation quality.

## 4. Training Dynamics Analysis

dLM training with the objective in Eq. 1 improves the masked denoising likelihood under noisy conditions, but how this improved likelihood estimation translates into downstream task accuracy and parallel token

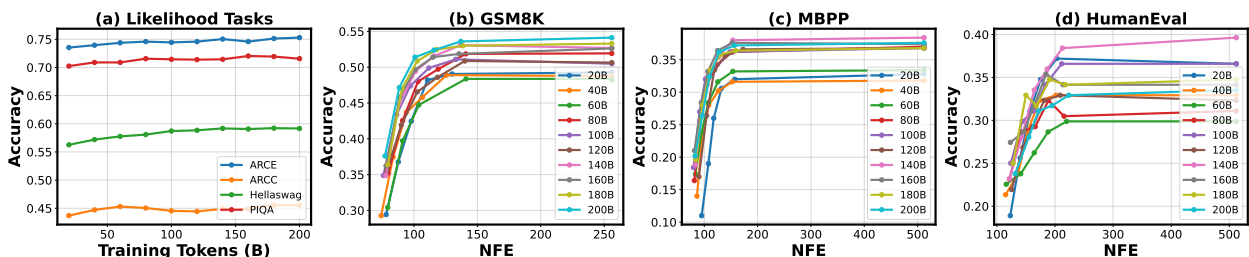


Figure 8 | (a) The accuracy evolution on likelihood tasks during training. (b–d) The accuracy–NFE trade-offs across different generation tasks for models trained with varying token budgets.

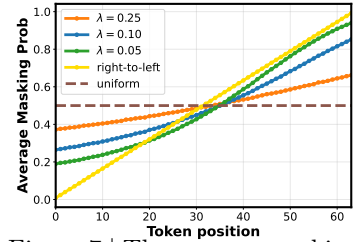


Figure 7 | The average masking probability of each token position within a block.

generation ability remains unclear. In this section, we study the training dynamics of dLMs by visualizing their task performance evolution during training.

**Setting.** We train Qwen2.5 1.5B for 200B tokens with the same setting as in Sec. 2.3, and evaluate on both generation and likelihood-based tasks, where accuracy is computed by estimating and selecting the largest likelihood among multiple choices. We visualize the accuracy evolution on likelihood tasks and the accuracy–efficiency (NFE) trade-off across different training token budgets in Fig. 8.

**Observations and analysis.** We observe that (1) with relatively low training cost (on the order of 10B tokens), dLMs converted from pretrained AR models can largely recover task accuracy. (2) Longer training with more iterations consistently improves likelihood estimation and yields higher accuracy on likelihood-based tasks. The average accuracy on generation tasks, without considering parallel token generation (i.e., the rightmost points of each curve in Fig. 8 (b–d)), also improves, though with fluctuations on certain tasks. (3) Improved likelihood estimation allows for more aggressive parallel token generation, as reflected in the enhanced accuracy–NFE trade-off with longer training. This indicates that stronger likelihood estimation produces more accurate and reliable confidence scores, thereby improving generation quality under confidence-based sampling.

This also indicates that parallel token generation ability is another dimension for evaluating a dLM’s performance: dLMs with comparable accuracy when denoising one token per step can exhibit notable accu-

racy gaps when performing aggressive parallel token generation.

**Takeaways:** dLMs’ ability to perform more aggressive parallel generation improves with better likelihood estimation, which can be induced by longer training on more tokens.

## 5. Efficient-DLM: A New Family of Efficient dLMs

Combining previous insights, we develop the Efficient-DLM family with three sizes (1.5B/4B/8B), continuously pretrained from Qwen2.5-1.5B (block size 16) and Qwen3 4B/Qwen3 8B (block size 64), respectively. Our Efficient-DLM family integrates (1) the identified best attention pattern, i.e., block-wise attention with clean context and without token shift, as analyzed in Sec. 2, and (2) the position-dependent token masking with  $\lambda = 0.1$  proposed in Sec. 3. Motivated by the training dynamics discussed in Sec. 4, we train for longer (300B tokens for Efficient-DLM 1.5B/4B and 500B tokens for Efficient-DLM 8B) using a mixed dataset comprising [12, 13, 14], adopting an initial learning rate of  $1e-5$  with cosine decay and the AdamW optimizer. All models are trained on 128 NVIDIA H100 GPUs.

### 5.1. Benchmark with SOTA AR and dLMs

We benchmark our Efficient-DLM against SOTA AR LMs (Qwen3 [16], Qwen2.5 [11], Llama3.2 [21], SmoLLM2 [22]) and SOTA dLMs (LLaDA [3] and Dream [4]) in Tab. 3. The benchmark covers 12 tasks, including math (GSM8k, Minerva Math), coding (HumanEval, HumanEval Plus, MBPP, MBPP Plus), factual knowledge (MMLU), and commonsense reasoning (ARCC, ARCE, Hellaswag, PIQA, Winogrande), as well as throughput measured on an NVIDIA H100 GPU with a batch size of 1. For each instance of Efficient-DLM, we report results under different parallel decoding settings with different tokens per forward, controlled by the confidence thresholds during confidence-based sampling [17]. More detailed settings are provided in Appendix A.

**Observations.** As shown in Tab. 3, we observe that (1) compared to SOTA dLMs, our Efficient-DLM achieves both higher accuracy and efficiency. For example, Efficient-DLM 8B delivers 5.35% higher average accuracy with  $4.50\times$  throughput over Dream 7B, benefiting from the block-wise attention design over fully bidirectional modeling in continuous pre-training (see Sec. 2.2). (2) Compared to SOTA AR LMs, our Efficient-DLM attains better accuracy-throughput trade-offs. For instance, Efficient-DLM 8B/4B achieves  $2.77\times/1.83\times$  throughput with  $+2.68\%/+7.79\%$  accuracy over Qwen3 4B/1.7B, re-

Table 3 | Benchmarking against SOTA AR models and dLMs on 12 tasks spanning coding, math, factual knowledge, and commonsense reasoning (CR), reporting average accuracy per category. TPF denotes tokens per forward, and TPS refers to throughput measured on an NVIDIA H100 GPU with a batch size of 1. Detailed per-task accuracy is provided in Appendix B.

Type	Model	TPF	TPS (tok/sec)	Coding	Math	MMLU	CR	Avg.
AR	Llama3.2 1B	1.00	143.91	24.45	4.98	30.98	60.62	34.24
	SmolLM2 1.7B	1.00	112.84	21.06	30.97	49.99	68.44	40.14
	Qwen2.5 0.5B	1.00	99.93	31.90	25.97	47.65	55.31	41.98
	Qwen2.5 1.5B	1.00	73.03	42.17	46.98	60.96	66.00	54.47
dLM	Efficient-DLM 1.5B	1.00	68.52	42.33	42.60	57.63	62.58	52.09
		2.33	158.89	42.33	42.28	57.63	62.58	52.04
		2.69	184.48	41.79	41.80	57.63	62.58	51.77
AR	Qwen3 1.7B	1.00	71.59	54.22	54.15	62.53	64.99	59.39
	Qwen3 4B	1.00	47.13	63.85	66.27	73.19	70.91	67.97
	Qwen3 8B	1.00	42.51	68.45	69.87	76.93	73.71	71.58
dLM	LLaDA 8B	1.00	25.04	38.10	49.13	65.86	68.50	54.92
	Dream 7B	1.00	28.11	58.92	58.39	67.00	72.83	65.30
dLM	Efficient-DLM 4B	1.00	44.13	62.08	67.98	71.80	70.87	67.54
		2.52	119.33	61.37	68.56	71.80	70.87	67.39
		3.01	130.24	60.96	68.10	71.80	70.87	67.18
dLM	Efficient-DLM 8B	1.00	39.99	67.36	69.22	77.22	74.88	71.62
		2.74	109.78	65.64	68.52	77.22	74.88	70.93
		3.27	130.71	64.95	68.21	77.22	74.88	70.65

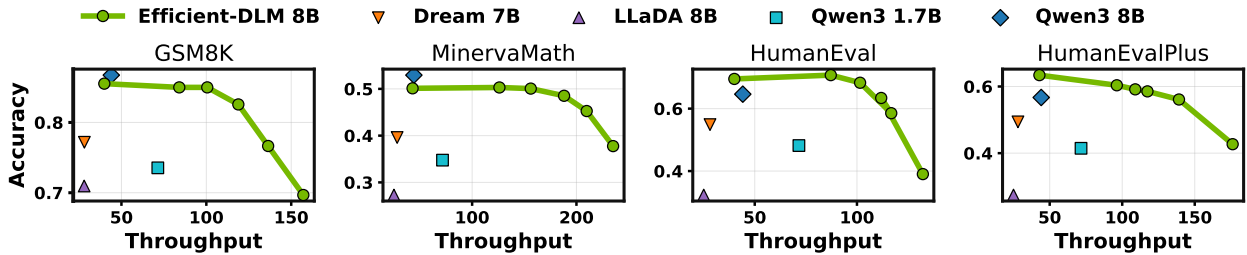


Figure 9 | Visualizing the accuracy-throughput trade-off of different models across different generation tasks.

spectively. More benchmarks with SOTA dLMs plus Fast-dLLM [17] are in Appendix C.

## 5.2. One-for-All Flexibility: Adaptive Accuracy–Efficiency Trade-offs

Beyond efficiency, another advantage of dLMs is their one-for-all flexibility: a single dLM can balance accuracy and throughput to suit different deployment scenarios. This is achieved by controlling parallel token generation via a confidence threshold [17]. Fig. 9 shows the one-for-all flexibility of our Efficient-DLM 8B across four math and coding tasks. We observe that a single Efficient-DLM 8B achieves better accuracy–throughput frontiers than the AR Qwen3 family from 1.7B to 8B, demonstrating its promise for one-for-all deployment. Throughput results under large batch sizes are provided in Appendix C.

## 5.3. Advantages of dLMs in Text Embedding

We further highlight that, thanks to their ability for bidirectional modeling, dLMs are more promising than AR models for tasks requiring high-quality text embeddings. To demonstrate this, we evaluate our Efficient-DLM against AR Qwen models on text embedding tasks, benchmarking 15 datasets from the MTEB benchmark [23] across six categories, following the ablation setup from LLM2Vec [24]. As shown in Tab. 4, we observe a clear advantage of dLMs: at the 1.5B and 4B scales, Efficient-DLM outperforms AR Qwen models of the same sizes by 7.71% and 9.91% on average, respectively. These results also highlight the broader promise of dLMs for other sequence modeling tasks that require bidirectional information.

Table 4 | Comparing our Efficient-DLM and AR Qwen models on text embedding tasks [23].

Model	Retr.	Ranking	Clust.	Pair Class.	Class.	STS	Avg.
Qwen2.5 1.5B	20.69	40.01	21.42	24.59	31.22	39.33	29.54
Efficient-DLM 1.5B	18.67	43.67	23.58	56.76	31.70	49.14	37.25
Qwen3 4B	19.46	39.90	21.77	33.94	29.13	40.56	30.79
Efficient-DLM 4B	20.17	45.05	23.91	65.59	42.22	47.27	40.70

## 5.4. Ablation Study of Different Components

We have analyzed and demonstrated the impact of each component of Efficient-DLM in Sec. 2–4. We also

summarize their impact in Tab. 5, which performs AR-to-dLM conversion on top of Qwen3 4B. We start from the baseline setting (Dream’s bidirectional modeling with 25B training tokens) and progressively add each component, reporting accuracy across six math and coding tasks. As observed in Tab. 5, proper attention patterns (with appropriate block-size selection), removing token shift, adding position-dependent token masking, and longer training all contribute to successful AR-to-dLM conversion.

## 6. Related Work

**Diffusion language models.** To overcome the token-by-token decoding nature of AR LMs, diffusion LMs, both continuous [25, 26, 27] and discrete [28, 1, 2, 29, 18], have been proposed to perform non-AR decoding and thus enable parallel token generation. Among them, masked dLMs [1, 2, 3, 4] have been successfully scaled up (e.g., LLaDA [3] and Dream [4]). Follow-up work has further explored alternative dLM paradigms [6, 7, 9], and scaled them to larger generalists [30] or domain-specific specialists such as coding agents [31, 32, 33, 34], explored dedicated reinforcement learning schemes [35, 36], and extended them to more modalities [37, 38]. Compared to AR LMs, diffusion LMs have been demonstrated to be better learners under data-constrained settings [39] and show improved performance in planning [4] and text embedding [40].

**Diffusion language model acceleration.** Despite the acceleration potential of large dLMs [3, 4], the gap between bidirectional attention and KV caching, along with the one-token-per-step denoising process, limits their achievable speed-up. To address these challenges, dedicated caching strategies for dLMs [41, 42, 17] have been developed to reuse computations and approximate bidirectional attention. In addition, to realize the potential of parallel token generation, confidence-based sampling [17], guidance from AR models [43], and adaptive decoding with certainty and positional priors [44] have been proposed. Beyond these training-free methods, [10, 4] propose initializing dLMs from AR models with token shifts to accelerate dLM training. Block Diffusion [5] combines AR and diffusion by performing block-wise AR and

Table 5 | Ablation study of different AR-to-dLM components on Qwen3 4B by progressively adding each component on top of the baseline setting (Dream’s bidirectional modeling with 25B training tokens).

Setting	HumanEval	HumanEval Plus	MBPP	MBPP Plus	GSM8K	Minerva Math	Avg
Bidirectional (Dream’s setting)	39.02	32.32	39.60	50.00	67.40	39.17	44.59
+ Block-wise w/ clean context	53.66	50.36	55.60	69.70	78.39	46.33	59.01
+ Remove token shift	56.10	51.22	54.60	69.84	82.87	47.02	60.27
+ Position-dependent masking	60.37	54.27	59.00	71.43	81.12	45.92	62.02
+ Scale to 300B tokens	60.98	56.71	60.00	70.63	86.43	49.54	64.05

in-block diffusion to support native KV caching, and concurrent works [17, 45, 46] also convert pretrained AR models or dLMs into block-wise dLMs.

## 7. Conclusion

This work systematically explores how to convert pretrained AR models into dLMs that achieve faster generation while retaining strong accuracy. By introducing a continuous pretraining scheme with a block-wise attention pattern, along with a position-dependent token masking strategy that narrows the training–test gap, we provide a principled framework for delivering dLMs with both strong accuracy and speed, resulting in the Efficient-DLM model family. Through comprehensive analyses of attention patterns, training dynamics, and other design choices, our findings offer actionable insights that we hope will guide the community toward building efficient and scalable dLMs. Beyond serving as a practical recipe for AR-to-dLM conversion, our results highlight the broader opportunity to rethink pretraining, masking, and decoding strategies for dLMs in order to realize their promise as alternatives to AR models.

## References

- [1] Zhengfu He, Tianxiang Sun, Kuanning Wang, Xuanjing Huang, and Xipeng Qiu. Diffusionbert: Improving generative masked language models with diffusion models. *arXiv preprint arXiv:2211.15029*, 2022.
- [2] Subham Sahoo, Marianne Arriola, Yair Schiff, Aaron Gokaslan, Edgar Marroquin, Justin Chiu, Alexander Rush, and Volodymyr Kuleshov. Simple and effective masked diffusion language models. *Advances in Neural Information Processing Systems*, 37:130136–130184, 2024.
- [3] 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. *arXiv preprint arXiv:2502.09992*, 2025.
- [4] Jiacheng Ye, Zhihui Xie, Lin Zheng, Jiahui Gao, Zirui Wu, Xin Jiang, Zhenguo Li, and Lingpeng Kong. Dream 7b: Diffusion large language models. *arXiv preprint arXiv:2508.15487*, 2025.
- [5] 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. *arXiv preprint arXiv:2503.09573*, 2025.
- [6] Subham Sekhar Sahoo, Zhihan Yang, Yash Akhauri, Johnna Liu, Deepansha Singh, Zhoujun Cheng, Zhengzhong Liu, Eric Xing, John Thickstun, and Arash Vahdat. Esoteric language models. *arXiv preprint arXiv:2506.01928*, 2025.
- [7] Subham Sekhar Sahoo, Justin Deschenaux, Aaron Gokaslan, Guanghan Wang, Justin Chiu, and Volodymyr Kuleshov. The diffusion duality. *arXiv preprint arXiv:2506.10892*, 2025.
- [8] Shen Nie, Fengqi Zhu, Chao Du, Tianyu Pang, Qian Liu, Guangtao Zeng, Min Lin, and Chongxuan Li. Scaling up masked diffusion models on text. *arXiv preprint arXiv:2410.18514*, 2024.
- [9] Shuchen Xue, Tianyu Xie, Tianyang Hu, Zijin Feng, Jiacheng Sun, Kenji Kawaguchi, Zhenguo Li, and Zhi-Ming Ma. Any-order gpt as masked diffusion model: Decoupling formulation and architecture. *arXiv preprint arXiv:2506.19935*, 2025.
- [10] 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. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [11] Qwen Team. Qwen2.5: A party of foundation models, September 2024.
- [12] NVIDIA Nemotron Nano. Efficient hybrid mamba-transformer reasoning model. *arXiv preprint arXiv:2508.14444*, 2025.
- [13] Fan Zhou, Zengzhi Wang, Nikhil Ranjan, Zhoujun Cheng, Liping Tang, Guowei He, Zhengzhong Liu, and Eric P. Xing. Megamath: Pushing the limits of open math corpora. *arXiv preprint arXiv:2504.02807*, 2025. Preprint.
- [14] Kazuki Fujii, Yukito Tajima, Sakae Mizuki, Hinari Shimada, Taihei Shiotani, Koshiro Saito, Masanari Ohi, Masaki Kawamura, Taishi Nakamura, Takumi

- Okamoto, Shigeki Ishida, Kakeru Hattori, Youmi Ma, Hiroya Takamura, Rio Yokota, and Naoaki Okazaki. Rewriting pre-training data boosts llm performance in math and code, 2025.
- [15] Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model evaluation harness, 07 2024.
- [16] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- [17] 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. *arXiv preprint arXiv:2505.22618*, 2025.
- [18] Jingyang Ou, Shen Nie, Kaiwen Xue, Fengqi Zhu, Jiacheng Sun, Zhenguo Li, and Chongxuan Li. Your absorbing discrete diffusion secretly models the conditional distributions of clean data. *arXiv preprint arXiv:2406.03736*, 2024.
- [19] Iris AM Huijben, Wouter Kool, Max B Paulus, and Ruud JG Van Sloun. A review of the gumbel-max trick and its extensions for discrete stochasticity in machine learning. *IEEE transactions on pattern analysis and machine intelligence*, 45(2):1353–1371, 2022.
- [20] Tong Wu, Zhihao Fan, Xiao Liu, Hai-Tao Zheng, Yeyun Gong, Jian Jiao, Juntao Li, Jian Guo, Nan Duan, Weizhu Chen, et al. Ar-diffusion: Autoregressive diffusion model for text generation. *Advances in Neural Information Processing Systems*, 36:39957–39974, 2023.
- [21] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [22] Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Penedo, Lewis Tunstall, Andrés Marafioti, Hynek Kydlicek, Agustín Piqueres Lajarín, Vaibhav Srivastav, Joshua Lochner, Caleb Fahlgrén, Xuan-Son Nguyen, Clémentine Fourrier, Ben Burtenshaw, Hugo Larcher, Haojun Zhao, Cyril Zakka, Mathieu Morlon, Colin Raffel, Leandro von Werra, and Thomas Wolf. Smollm2: When smol goes big – data-centric training of a small language model, 2025.
- [23] Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. Mteb: Massive text embedding benchmark. *arXiv preprint arXiv:2210.07316*, 2022.
- [24] Parishad BehnamGhader, Vaibhav Adlakha, Marius Mosbach, Dzmitry Bahdanau, Nicolas Chapados, and Siva Reddy. LLM2Vec: Large language models are secretly powerful text encoders. In *First Conference on Language Modeling*, 2024.
- [25] Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-lm improves controllable text generation. *Advances in neural information processing systems*, 35:4328–4343, 2022.
- [26] Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, and LingPeng Kong. Diffuseq: Sequence to sequence text generation with diffusion models. *arXiv preprint arXiv:2210.08933*, 2022.
- [27] Xiaochuang Han, Sachin Kumar, and Yulia Tsvetkov. Ssd-lm: Semi-autoregressive simplex-based diffusion language model for text generation and modular control. *arXiv preprint arXiv:2210.17432*, 2022.
- [28] Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in neural information processing systems*, 34:17981–17993, 2021.
- [29] Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion modeling by estimating the ratios of the data distribution. *arXiv preprint arXiv:2310.16834*, 2023.
- [30] Google DeepMind. Gemini diffusion, 2025. Model page: state-of-the-art, experimental text diffusion model.
- [31] Samar Khanna, Siddhant Kharbanda, Shufan Li, Harshit Varma, Eric Wang, Sawyer Birnbaum, Ziyang Luo, Yanis Miraoui, Akash Palrecha, Stefano Ermon, et al. Mercury: Ultra-fast language models based on diffusion. *arXiv preprint arXiv:2506.17298*, 2025.
- [32] 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.
- [33] 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.
- [34] Zihui 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.
- [35] Siyan Zhao, Devaansh Gupta, Qinqing Zheng, and Aditya Grover. d1: Scaling reasoning in diffusion large language models via reinforcement learning. *arXiv preprint arXiv:2504.12216*, 2025.

- [36] 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.
- [37] 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.
- [38] Ling Yang, Ye Tian, Bowen Li, Xinchun Zhang, Ke Shen, Yunhai Tong, and Mengdi Wang. Mmada: Multimodal large diffusion language models. *arXiv preprint arXiv:2505.15809*, 2025.
- [39] Mihir Prabhudesai, Mengning Wu, Amir Zadeh, Kate-rina Fragkiadaki, and Deepak Pathak. Diffusion beats autoregressive in data-constrained settings. *arXiv preprint arXiv:2507.15857*, 2025.
- [40] Siyue Zhang, Yilun Zhao, Liyuan Geng, Arman Cohen, Anh Tuan Luu, and Chen Zhao. Diffusion vs. autoregressive language models: A text embedding perspective. *arXiv preprint arXiv:2505.15045*, 2025.
- [41] 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.
- [42] Xinyin Ma, Runpeng Yu, Gongfan Fang, and Xinchao Wang. dkv-cache: The cache for diffusion language models. *arXiv preprint arXiv:2505.15781*, 2025.
- [43] Daniel Israel, Guy Van den Broeck, and Aditya Grover. Accelerating diffusion llms via adaptive parallel decoding. *arXiv preprint arXiv:2506.00413*, 2025.
- [44] 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.
- [45] Shuang Cheng, Yihan Bian, Dawei Liu, Linfeng Zhang, Qian Yao, Zhongbo Tian, Wenhai Wang, Qipeng Guo, Kai Chen, Biqing Qi, et al. Sdar: A synergistic diffusion-autoregression paradigm for scalable sequence generation. *arXiv preprint arXiv:2510.06303*, 2025.
- [46] 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.
- [47] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.

## A. Detailed Experimental Settings

**Evaluation settings.** For all evaluations of our Efficient-DLM in Sec. 5.1, we follow the best practice from Sec. 2.3 and use evaluation block sizes of 16 for Efficient-DLM 1.5B and 32 for Efficient-DLM 4B/8B, respectively. We use lm-evaluation-harness [15] to evaluate AR baselines (Qwen3 [16], Qwen2.5 [11], Llama3.2 [21], SmolLM2 [22]); for dLMs (LLaDA [3] and Dream [4]), we adopt their official evaluation code. We benchmark 12 tasks covering math (GSM8k, Minerva Math), coding (HumanEval, HumanEval Plus, MBPP, MBPP Plus), factual knowledge (MMLU), and commonsense reasoning (ARCC, ARCE, Hellaswag, PIQA, Winogrande). Following Dream [4], we use 8-shot, 4-shot, 0-shot, 0-shot, 3-shot, and 3-shot settings for GSM8k, Minerva Math, HumanEval, HumanEval Plus, MBPP, and MBPP Plus, respectively. The maximum number of generated tokens is set to 512 for all tasks, except GSM8k, which uses 256 as in [4].

**Parallel decoding settings.** Following [17], we set a confidence threshold and decode all tokens that exceed the threshold at each denoising step to enable parallel decoding. Tokens per forward (TPF) and generation throughput (tok/sec), reported in Tab. 3, are averaged over all six generation tasks.

**Text embedding evaluation settings.** Following the ablation setting of LLM2Vec [24], we evaluate the models on six categories of tasks from MTEB [23], including retrieval (SciFact, ArguAna, NFCorpus), reranking (StackOverflowDupQuestions, SciDocsRR), clustering (BiorxivClusteringS2S, MedrxivClusteringS2S, TwentyNewsgroupsClustering), pair classification (SprintDuplicateQuestions), classification (Banking77Classification, EmotionClassification, MassiveIntentClassification), and semantic textual similarity (STS17, SICK-R, STSBenchmark). In total, the evaluation covers 15 datasets.

To obtain sequence embeddings, we apply mean pooling over the last layer’s hidden states across all tokens. For Qwen models, we use a causal attention mask, as switching to a bidirectional mask consistently degraded performance. For our Efficient-DLM models, we instead employ a bidirectional attention mask. All experiments are conducted in a zero-shot setting, directly using the pretrained weights of Qwen and Efficient-DLM without any fine-tuning.

## B. Per-task Accuracy Achieved by Efficient-DLM and Baselines

We provide the per-task accuracy achieved by Efficient-DLM and SOTA AR/dLMs in Tab. 6 as

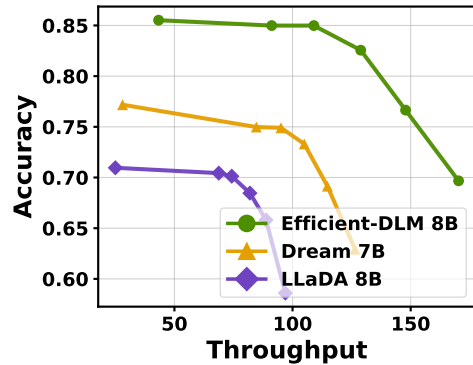


Figure 10 | Comparing the accuracy-throughput trade-off with Dream/LLaDA plus Fast-dLLM.

a complement to Tab. 3 of our main paper.

## C. More Benchmarks with SOTA AR LMs and dLMs

**Benchmark with SOTA dLMs plus Fast-dLLM acceleration.** Existing public dLMs, LLaDA [3] and Dream [4], both equipped with fully bidirectional attention, can be accelerated by Fast-dLLM [17] through partial KV caching and parallel decoding. We benchmark our Efficient-DLM 8B against Dream and LLaDA enhanced with Fast-dLLM’s dual cache and parallel decoding, using different confidence thresholds to control the accuracy–throughput trade-off. As shown in Fig. 10, our Efficient-DLM 8B consistently achieves a better accuracy–throughput trade-off than Dream and LLaDA on GSM8K, demonstrating that Efficient-DLM constitutes a stronger dLM family.

**The accuracy–throughput trade-off with larger inference batch sizes.** As a complement to Sec. 5.2, we further visualize the trade-off between accuracy and throughput under different batch sizes for multiple models on the GSM8K dataset. As shown in Fig. 11, we observe that (1) our Efficient-DLM 8B consistently improves the accuracy–efficiency trade-off compared to both AR models (Qwen3 1.7B–8B) and dLMs, up to a batch size of 16; and (2) the efficiency benefits of dLMs over AR models are more pronounced at small batch sizes, which correspond to more memory-bounded scenarios, and these benefits begin to diminish at larger batch sizes, e.g., Efficient-DLM 8B falls behind Qwen3 1.7B in throughput at a batch size of 32.

This set of experiments highlights the current limitations of dLMs in large-batch serving scenarios. Potential workarounds include adaptive block sizes, improved parallel sampling schemes, and combining dLMs with linear attention to enhance large-batch efficiency, which we leave for future work.

Table 6 | Benchmarking against SOTA AR models and dLMs on 12 tasks spanning coding, math, factual knowledge, and commonsense reasoning. TPF denotes tokens per forward, and TPS refers to throughput measured on an NVIDIA H100 GPU with a batch size of 1. This table is a complement to Tab. 3.

Type	Model	TPF	Throughput (t/s)	Coding				Math		Factual		Commonsense Reasoning			Avg	
				HumanEval	HumanEval Plus	MBPP	MBPP Plus	GSM8K	Minerva Math	MMLU	ARC-E	ARC-C	Hellaswag	PIQA		Winogrande
AR	Llama3.2 1B	1.00	138.72	17.68	14.63	26.60	38.89	5.69	4.28	30.98	65.28	36.35	63.76	74.43	63.30	36.82
	Qwen2.5 0.5B	1.00	95.40	27.44	25.61	29.60	44.97	37.45	14.48	47.65	64.77	31.83	52.25	70.02	57.70	41.98
	Qwen2.5 1.5B	1.00	73.03	35.98	29.88	43.60	59.23	64.97	28.98	60.96	75.17	45.05	67.90	76.12	65.75	54.47
dLM (Ours)	Efficient-DLM 1.5B	2.33	158.89	41.46	35.37	38.80	53.70	57.92	26.64	57.63	75.00	45.31	59.43	71.44	61.72	52.04
		2.69	184.48	41.46	35.98	37.60	52.12	58.07	25.52	57.63	75.00	45.31	59.43	71.44	61.72	51.77
AR	Qwen3 1.7B	1.00	71.59	48.17	41.46	55.80	71.43	73.54	34.76	62.53	73.32	44.62	66.43	75.63	64.96	59.39
	Qwen3 4B	1.00	44.99	57.32	50.61	66.80	83.69	85.44	47.10	73.19	79.12	51.62	73.66	78.07	72.06	67.97
	Qwen3 8B	1.00	42.51	64.63	56.71	69.40	83.07	86.73	52.94	76.93	81.90	53.16	78.59	79.22	75.69	71.58
dLM	LLaDA 8B	1.00	25.04	32.32	27.44	40.80	51.85	70.96	27.30	65.86	73.78	49.15	71.05	73.88	74.66	54.92
		1.00	28.11	54.88	49.39	56.80	74.60	77.18	39.60	67.00	82.20	59.13	73.73	75.52	73.56	65.30
dLM (Ours)	Efficient-DLM 4B	2.52	119.33	61.59	56.71	57.60	69.58	87.57	49.56	71.80	81.52	55.80	69.02	75.46	72.53	67.39
		3.01	130.24	60.98	56.10	57.20	69.58	87.19	49.02	71.80	81.52	55.80	69.02	75.46	72.53	67.18
dLM (Ours)	Efficient-DLM 8B	1.00	40.00	68.29	63.41	61.00	76.72	88.32	50.12	77.22	84.89	61.86	72.53	77.53	77.58	71.62
		2.74	109.78	64.02	59.15	63.20	76.19	86.73	50.30	77.22	84.89	61.86	72.53	77.53	77.58	70.93
		3.27	130.71	63.41	58.54	62.20	75.66	86.35	50.06	77.22	84.89	61.86	72.53	77.53	77.58	70.65

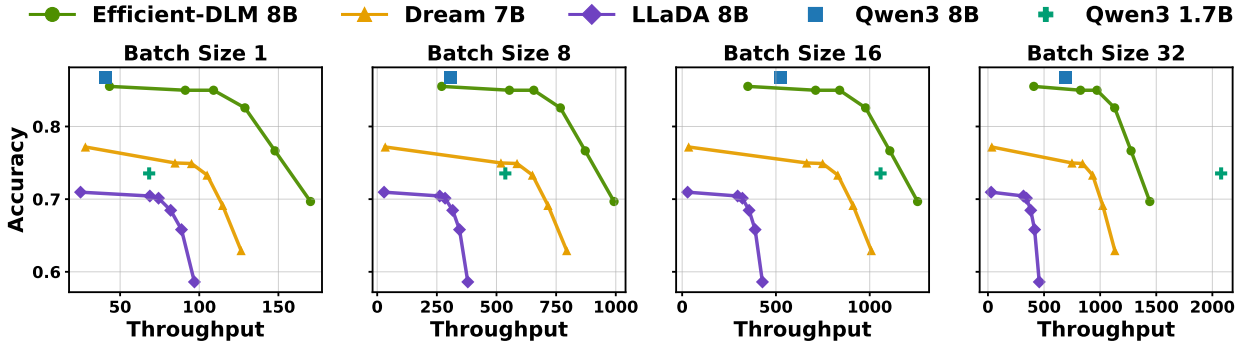


Figure 11 | Visualizing the accuracy-throughput trade-off under different batch sizes on the GSM8K dataset.

## D. The Impact of Initial LR

When initializing from pretrained AR models, the learning rate for continuous pretraining is a key hyperparameter, as it controls the speed of weight changes that affect both the preservation of the pretrained models’ abilities and the adaptation to dLMs’ new attention patterns. We perform an ablation study on Qwen3 4B trained for 25B tokens with different initial learning rates using a cosine learning rate schedule.

**Observations and analysis.** As shown in Tab. 7, we find that there exists a sweet-spot learning rate setting, e.g.,  $1e-5$  in our case, that balances both aspects mentioned above. Intuitively, overly large learning rates cause greater weight drifts and degrade the pretrained models’ original abilities, while overly small learning rates cannot effectively adapt to the new attention pattern. We also note that for any design factors in continuously training a pretrained model into dLMs, these two aspects should be carefully balanced to achieve decent final accuracy. Based on this set of experiments, we adopt  $1e-5$  as the default initial learning rate throughout the main manuscript.

## E. Loss Distributions across Tokens

To study the difference in loss distributions across token positions between AR models and dLMs, we

visualize the training loss defined in Eq. 1 for diffusion Qwen2.5 1.5B trained with a block size of 16, alongside the AR Qwen2.5 1.5B trained with an AR loss.

As shown in Fig. 12, which shows the average loss of the first 256 tokens in training sequences, we observe that (1) in AR models, the initial tokens incur higher loss due to the lack of context, while the loss of later tokens becomes more uniform; and (2) in dLMs, the loss follows a periodic pattern aligned with block boundaries, where later tokens within each block show higher loss due to limited clean context, consistent with Fig. 6 (c). In addition, similar to AR models, the initial tokens of the entire sequence also experience higher loss from insufficient context.

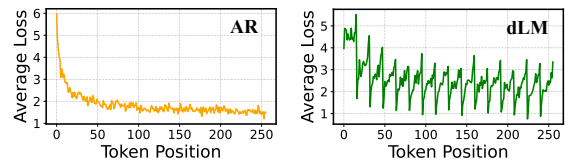


Figure 12 | Visualizing the loss distributions over token positions of AR models and dLMs.

## F. AR-to-dLM Conversion via Parameter-Efficient Tuning

Motivated by the relatively small weight changes observed in Sec. 2.2, we investigate whether parameter-

Table 7 | Results of continuous pretraining with different initial learning rates on Qwen3 4B for 25B tokens.

Init LR	HumanEval	HumanEval Plus	MBPP	MBPP Plus	GSM8K	Minerva Math	Avg
1.00E-04	49.39	43.90	44.20	56.08	72.56	39.54	50.95
3.00E-05	54.27	49.39	52.40	67.46	77.48	38.56	56.59
1.00E-05	57.93	51.22	54.40	71.96	81.73	46.54	60.63
3.00E-06	56.10	50.61	54.60	67.99	83.93	47.44	60.11
1.00E-06	45.73	42.68	47.20	66.67	81.12	43.94	54.56

Table 8 | Comparison of different dLM training schemes on Qwen2.5 1.5B. This table extends Tab. 1, with Rows (h) and (i) added to present the LoRA tuning results.

Row ID	Attn Pattern	Clean Context	Token Shift	KV Cache	LoRA Rank	Human -Eval	Human -Eval Plus	MBPP	MBPP Plus	GSM8K	Minerva Math	Avg
a	AR	-	✓	✓	-	36.59	29.88	43.6	59.52	54.74	26.40	41.79
b	Bidirectional	-	✓	✗	-	15.85	12.20	16.2	24.34	28.96	11.08	18.10
c	Bidirectional	-	✗	✗	-	19.51	15.24	17.2	24.34	28.20	11.22	19.29
d	Block-wise	✗	✓	✓	-	31.10	25.61	23.6	36.77	38.44	13.88	28.23
e	Block-wise (2×)	✗	✓	✓	-	26.22	22.56	26.0	42.33	36.69	12.56	27.73
f	Block-wise	✓	✓	✓	-	38.41	33.54	33.0	48.68	51.48	21.04	37.69
g	Block-wise	✓	✗	✓	-	39.02	34.76	34.0	48.15	52.99	21.56	38.41
h	Block-wise (LoRA)	✓	✗	✓	16	30.49	25.61	20.60	30.95	43.82	16.08	27.93
i	Block-wise (LoRA)	✓	✗	✓	64	28.66	25.61	24.40	40.21	48.14	17.64	30.78

efficient tuning can effectively convert pretrained AR models into dLMs. To this end, we apply Low-Rank Adaptation (LoRA) [47] to all linear layers in attention/FFN modules, combined with the best attention pattern identified in Sec. 2.2, i.e., block-wise attention conditioning on clean context without token shift. All other parameters are frozen, except for the embedding layer, normalization operators, and the final model head, which we find must remain trainable for effective adaptation.

**Observations and analysis.** We extend Tab. 1 into Tab. 8 by including LoRA tuning results with two different ranks in Rows (h) and (i). We observe that LoRA tuning achieves reasonably good performance for AR-to-dLM conversion. Specifically, LoRA with rank 64 (Row i) surpasses the full-model training results of fully bidirectional attention and block-wise attention without clean context, while remaining 7.63% behind the full-model training results of the best scheme, i.e., block-wise attention with clean context. These results indicate that (1) with proper training schemes, even parameter-efficient tuning can yield competitive dLMs, and (2) full-model training remains necessary to obtain strong dLMs.