

Causal Autoregressive Diffusion Language Model

Junhao Ruan^{1,2} Bei Li² Yongjing Yin² Pengcheng Huang¹ Xin Chen² Jingang Wang² Xunliang Cai²
Tong Xiao^{1,3} Jingbo Zhu^{1,3}

Abstract

In this work, we propose Causal Autoregressive Diffusion (CARD), a novel framework that unifies the training efficiency of ARMs with the high-throughput inference of diffusion models. CARD reformulates the diffusion process within a strictly causal attention mask, enabling dense, per-token supervision in a single forward pass. To address the optimization instability of causal diffusion, we introduce a soft-tailed masking schema to preserve local context and a context-aware reweighting mechanism derived from signal-to-noise principles. This design enables dynamic parallel decoding, where the model leverages KV-caching to adaptively generate variable-length token sequences based on confidence. Empirically, CARD outperforms existing discrete diffusion baselines while reducing training latency by $3 \times$ compared to block diffusion methods. Our results demonstrate that CARD achieves ARM-level data efficiency while unlocking the latency benefits of parallel generation, establishing a robust paradigm for next-generation efficient LLMs.

1. Introduction

Causal Autoregressive Models (ARMs) currently serve as the dominant paradigm for training Large Language Models (LLMs), owing to their stable training dynamics and predictable scaling laws. However, as model parameters and test-time compute requirements grow, the sequential nature of autoregressive decoding has emerged as a critical bottleneck. This inefficiency has sparked renewed interest in Text Diffusion Models, which offer theoretical advantages including parallel inference (Austin et al., 2021), iterative refinement (Wang et al., 2025), and potentially higher data

¹School of Computer Science and Engineering, Northeastern University, Shenyang, China ²Meituan Inc. ³NiuTrans Research, Shenyang, China. Correspondence to: Bei Li <libei17@meituan.com>, Tong Xiao <xiao-tong@mail.neu.edu.cn>.

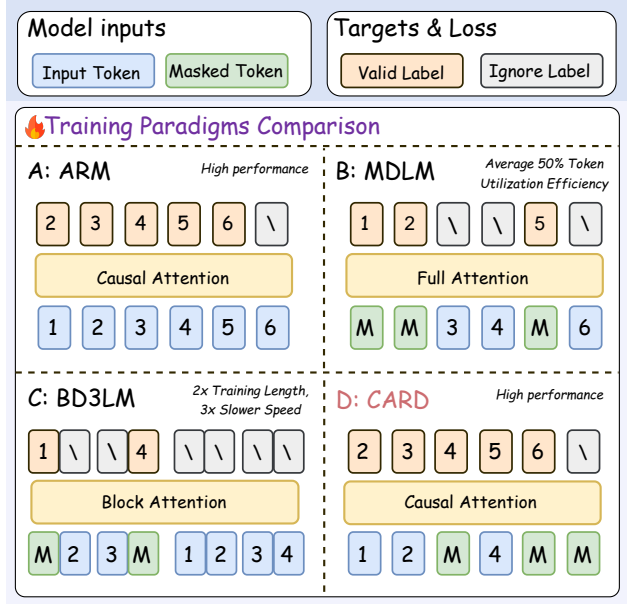


Figure 1. Comparison of training paradigms. Current diffusion methods like MDLM and BD3LM are inefficient compared to ARM; MDLM reaches only 50% of ARM’s expected efficiency, while BD3LM relies on complex masking and sequence duplication. CARD overcomes these issues by using causal diffusion, maintaining the same high efficiency as ARM while achieving better performance.

modeling capacity (Ni et al., 2025).

Early attempts at discrete diffusion faced significant hurdles due to complex training objectives involving variational bounds and numerical instabilities caused by noise sampling (Austin et al., 2021). A turning point occurred with the introduction of Simplified Masked Discrete Diffusion Models (MDLM) (Shi et al., 2024; Sahoo et al., 2024). By simplifying the diffusion process into a subspace assumption analogous to a randomized Masked Language Modeling (MLM) (Devlin et al., 2019) task, MDLM ushered in the era of scalable text diffusion (Nie et al., 2025a), enabling the training of modern LLM-scale diffusion models like LLaDA (Nie et al., 2025b) and Dream (Ye et al., 2025).

Despite these advancements, standard MDLMs face severe architectural constraints. As illustrated in the MDLM panel of Figure 1, its reliance on bidirectional (“Full”) attention

Autoregressive

Generation steps **High Quality, Arbitrary-length, KV caching, Not Parallelizable**

↓
Nuclear energy is the cornerstone
Nuclear energy is the cornerstone of
Nuclear energy is the cornerstone of french ...

Block Diffusion Model

Generation steps **Lower Quality, Arbitrary-length, KV caching, Constant Parallelism**

↓
The names chiseled onto
The names chiseled onto city tenement
The names chiseled onto city tenement building entrances

Mask Diffusion Language Model

Generation steps **Lower Quality, Fixed-length, No KV caching, Adaptive Parallelism**

↓
Helplessness confusion are
Helplessness and confusion are easily come to
Helplessness and confusion are words that easily come to mind

Causal Diffusion Language Model (Ours)

Generation steps **High Quality, Arbitrary-length, KV caching, Adaptive Parallelism**

↓
The names chiseled
The names chiseled onto city
The names chiseled onto city tenement building entrances

Figure 2. Inference comparison of the four paradigms. CARD achieves high-quality results similar to ARM. With KV cache support, friendly operators, and parallel generation, it offers faster throughput than earlier methods. In particular, our inference parallelism is flexible, unlike BD3LM which is tied to the fixed block size used during training.

prevents the utilization of Key-Value (KV) caching. Consequently, inference speed often falls behind ARMs in practical scenarios (Wu et al., 2025). Furthermore, the arbitrary dependency order in training can lead to ineffective learning pathways (Kim et al., 2025), and the architecture fundamentally lacks support for variable-length generation.

To address these limitations, recent works have proposed hybrid architectures such as Block Diffusion (e.g., BD3LM) (Arriola et al., 2025). These models (drawn in Figure 1) operate at a coarser granularity, applying causal attention between fixed-size blocks and bidirectional attention within them. However, it introduces significant computational overhead. The vectorization required for block-wise training necessitates complex attention masking and can increase memory consumption and training latency by factors of $2\times$ and $3\times$, respectively. Moreover, the rigid, fixed block size fails to adapt to the varying information density inherent in natural language, limiting dynamic parallelism.

In this work, we propose CARD, a framework that combines the training efficiency of ARMs with the parallel inference of diffusion models through a strictly causal formulation. For training, CARD employs a *shifted causal attention mechanism* where each position predicts its original token from the preceding noised context. This generates a dense diffusion loss for the entire sequence in a single forward pass, achieving 100% token utilization without the overhead of block vectorization. For inference, CARD’s causal structure enables KV-caching (Figure 2), allowing the model to append a variable number of [MASK] tokens to the prefix and decode them in parallel through iterative denoising. This dynamic strategy generates multiple tokens per step when confidence is high while falling back to sequential decoding when necessary.

We empirically validate CARD on 1B-parameter models trained on 300B tokens, benchmarking against state-of-the-art autoregressive and diffusion baselines. Our results demonstrate that CARD effectively bridges the gap between efficiency and performance:

- **Superior Performance:** CARD achieves an average zero-shot accuracy of **53.2%**, outperforming existing diffusion models (MDLM and BD3LM) by over **5.7 points** and matching the generation quality of ARMs. Notably, it achieves the lowest zero-shot perplexity on 6 out of 8 evaluated domains.
- **Training & Inference Efficiency:** By eliminating block-wise overhead, CARD reduces training latency by $3\times$ compared to Block Diffusion, matching the throughput of standard ARMs. During inference, our confidence-based decoding achieves **$1.7\times$ to $4.0\times$** wall-clock speedup with negligible quality degradation.
- **Data Potential:** Scaling analysis reveals that CARD possesses higher data efficiency than ARMs in data-constrained settings, continuing to improve performance through repeated training epochs where autoregressive baselines saturate.

2. Background

We review the evolution of text diffusion models and the specific discrete objective function that serves as the foundation for our work.

2.1. Evolution of Text Diffusion Models

Applying diffusion to the discrete domain of language has followed two primary trajectories: continuous embedding methods and discrete state-space models. Continuous approaches, such as Diffusion-LM (Li et al., 2022) and Dif-fuSeq (Gong et al., 2023), map discrete tokens to Gaussian latent spaces. The disconnect between the continuous diffusion process and the discrete nature of text leads to rounding errors during decoding, often resulting in lower generation performance compared to autoregressive baselines.

Discrete DDPM (D3PM) (Austin et al., 2021) addressed this by defining the corruption process directly on the vocabulary via transition matrices. While theoretically rigorous, D3PMs initially suffered from optimization instability and

inefficient inference. To mitigate this, SEDD (Lou et al., 2024) reformulated the objective using score entropy, aligning discrete diffusion closer to its continuous counterparts. However, SEDD relied on time-dependent probability ratios, which prevented step-skipping and slowed inference. RADD (Ou et al., 2025) later demonstrated that the explicit time dependency in the input was not strictly necessary for mathematical validity, enabling flexible sampling strategies.

A paradigm shift occurred with the introduction of MDLM (Sahoo et al., 2024) and MD4 (Shi et al., 2024). By isolating the absorbing state (masking) transition, these works reduced the complex variational bound to a simplified, randomized Masked Language Modeling (MLM) objective. This simplification significantly improved numerical stability and allowed for scaling laws to be established (Nie et al., 2025a), culminating in large-scale pre-trained models like LLaDA (Nie et al., 2025b).

Despite these successes, standard MDLMs utilize bidirectional attention, which prevents the use of KV caching and degrades inference speed for long sequences. BD3LM (Arriola et al., 2025) attempts to bridge this gap by segmenting sequences into fixed-size blocks with causal masking between them. While this restores some parallel generation capabilities, BD3LM imposes significant training overheads due to complex attention masks and input duplication. Semi-autoregressive architectures have been further explored in works like LLaDA2 (Bie et al., 2025), SDAR (Cheng et al., 2025).

The concurrent WeDLM (Liu et al., 2025) further specializes this by employing unidirectional attention within blocks; however, it generally adheres to the block diffusion paradigm where training operates at the block level rather than the token level which will also bring extra training cost. Distinctly, another concurrent work, C²DLM (Han et al., 2025), explores causality through the lens of semantic concepts rather than model architecture. It analyzes causal relationships within training data but retains a bidirectional backbone and the standard MDLM training objective.

2.2. Discrete Diffusion Formulation

Our method builds upon the absorbing state diffusion framework. Let \mathbf{x}_0 be a sequence of length L . D3PM optimizes the negative Variational Lower Bound (ELBO) over T steps, which decomposes into:

$$\begin{aligned} L_{\text{vb}} = & \underbrace{D_{\text{KL}}[q(\mathbf{x}_T|\mathbf{x}_0)||p(\mathbf{x}_T)]}_{L_T} - \underbrace{\mathbb{E}_q[\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]}_{L_0} \\ & + \sum_{t=2}^T \underbrace{\mathbb{E}_q[D_{\text{KL}}[q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)||p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)]]}_{L_{t-1}}. \end{aligned} \quad (1)$$

In practice, to improve training stability and sample quality,

D3PM often incorporates an auxiliary cross-entropy loss to directly predict \mathbf{x}_0 :

$$L_{\text{D3PM}} = L_{\text{vb}} + \lambda \mathbb{E}_{q(\mathbf{x}_0)} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [-\log \tilde{p}_\theta(\mathbf{x}_0|\mathbf{x}_t)], \quad (2)$$

where λ is a hyperparameter balancing the two terms.

This objective involves summing over the entire vocabulary for the posterior computation, making it computationally expensive. MDLM drastically simplifies this by employing a SUBS parameterization (Sahoo et al., 2024), where the model predicts \mathbf{x}_0 directly and unmasked tokens are carried over deterministically. The KL divergences collapse, and in the continuous-time limit ($T \rightarrow \infty$), the loss becomes a weighted MLM objective:

$$\mathcal{L}_{\text{MDLM}} = \mathbb{E}_{t \sim \mathcal{U}[0,1]} \left[w(t) \sum_{\ell \in \mathcal{M}_t} \log p_\theta(x^\ell | \mathbf{x}_t, t) \right]. \quad (3)$$

Here, the loss is computed only over the masked tokens \mathcal{M}_t . The weighting term $w(t) = \frac{\alpha'_t}{1-\alpha_t}$ is determined by the noise schedule α_t .

BD3LM (Arriola et al., 2025) extends this formulation to interpolate between autoregression and diffusion. By partitioning the sequence \mathbf{x} into B blocks, BD3LM defines an autoregressive distribution over blocks while performing discrete diffusion *within* each block. The objective applies the MDLM loss per block, conditioned on the clean history of previous blocks $\mathbf{x}^{<b}$:

$$\mathcal{L}_{\text{BD3LM}}(\mathbf{x}; \theta) = \sum_{b=1}^B \mathbb{E}_{t \sim \mathcal{U}[0,1]} \mathbb{E}_q [w(t) \log p_\theta(\mathbf{x}^b | \mathbf{x}_t^b, \mathbf{x}^{<b})], \quad (4)$$

where \mathbf{x}_t^b represents the noisy state of block at time t .

2.3. Absorbing State Diffusion Process

In this work, we focus on the specific discrete diffusion process that serves as our foundation. Unlike continuous diffusion models that operate on Gaussian noise, text diffusion models typically define a corruption process over a discrete vocabulary. We consider a continuous-time variable $t \in [0, 1]$, where $t = 0$ corresponds to the clean ground-truth sequence \mathbf{x}_0 , and $t = 1$ represents a fully masked sequence.

We utilize the absorbing state (masking) transition. For any token in the sequence at time t , the forward process determines whether it remains its original value or transitions to a special [MASK] token. This is governed by a noise schedule $\sigma(t)$. For a given t , each token is independently replaced by [MASK] with probability $P(x^t = [\text{MASK}] | x_0) = \sigma(t)$. Throughout this paper, we adopt a linear schedule where $\sigma(t) = t$.

This formulation allows us to bridge the gap between deterministic text and stochastic training. At any step t , the

Algorithm 1 CARD Training Framework

```

1: Input: Sequence  $\mathbf{x}_0$ , Model  $\theta$ 
2: Params: Tail factor  $\lambda$ , Base  $\beta$ , Decay  $p$ 
3: // 1. Noise Scheduling
4: Sample  $t \sim \mathcal{U}[0, 1]$ 
5: // 2. Soft Tail Masking
6:  $N = \max(1, \lfloor L \cdot t \rfloor)$ ,  $W = \min(L, \lfloor N \cdot \lambda \rfloor)$ 
7: Define tail window indices:  $\mathcal{I}_{\text{win}} = \{L - W + 1, \dots, L\}$  and
   sample a subset of indices  $\mathcal{M} \subset \mathcal{I}_{\text{win}}$  such that  $|\mathcal{M}| = N$ 
8: Initialize  $\mathbf{x}^t = \mathbf{x}_0$ 
9: for each  $n \in \mathcal{M}$  do
10:    $x_n^t \leftarrow [\text{MASK}]$ 
11: end for
12: // 3. Context-aware Reweighting
13: for  $n = 1$  to  $L$  do
14:    $C_n = \mathbb{I}[x_n^t \text{ is [MASK]}] \cdot (1 + \mathbb{I}[x_{n-1}^t \text{ is [MASK]}])$ 
15:    $S_n^{\text{local}} = \sum_{i=1}^n C_i \cdot (1 - p)^{(n-1-i)}$ 
16:    $w_n = (\beta + S_n^{\text{local}})^{-1}$ 
17: end for
18: // 4. Optimization
19:  $\mathcal{L}_{\text{CARD}} = \sum_{n=1}^L w_n \log p_{\theta}(x_0, n | \mathbf{x}_{<n}^t)$ 
20: Update  $\theta$  using  $\nabla_{\theta} \mathcal{L}_{\text{CARD}}$ 
    
```

model is presented with a partially corrupted version of the input, denoted as \mathbf{x}^t . The training objective is to learn a denoising function that recovers the original tokens \mathbf{x}_0 from these noisy observations. By sampling t uniformly during training, the model learns to handle varying levels of corruption, from simple text completion to complex generation from scratch.

3. The CARD Framework

We propose the Causal Autoregressive Diffusion (CARD) framework, the overall training procedure of which is summarized in Algorithm 1. CARD utilizes a continuous-time noise addition method to apply diffusion processes within a causal architecture. This approach allows the model to leverage the robustness of diffusion training while maintaining the efficiency of autoregressive generation.

3.1. Synthesizing Autoregression and Diffusion

The core philosophy of CARD is to unify the stable training dynamics of ARMs with the flexible generation capabilities of Diffusion Models. We achieve this synthesis via a *shifted causal attention mechanism*. Unlike standard ARMs that condition on a static, clean history to model $p(x_n | \mathbf{x}_{<n})$, CARD predicts the original token x_n conditioned on a corrupted prefix $\mathbf{x}_{<n}^t$ sampled from a continuous-time diffusion process. This architecture allows us to strictly maintain the triangular attention mask inherent to GPT-style models for computational efficiency, while simultaneously minimizing the expected reconstruction error across varying noise intensities. We formally define the resulting optimization objective, which aggregates the weighted log-likelihoods

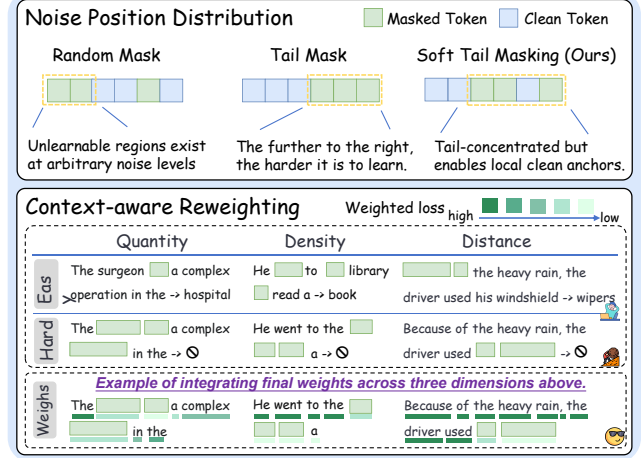


Figure 3. Soft Tail Masking concentrates noise at the sequence tail to resolve unlearnable regions in causal models via local clean anchors. (Bottom) Context-aware Reweighting adaptively down-weights the loss for high-ambiguity contexts similar to the diffusion ELBO principle, improving training stability.

across all token positions, as follows:

$$\mathcal{L}_{\text{CARD}} = \mathbb{E}_{t \sim \mathcal{U}[0,1], \mathbf{x}^t \sim q(\mathbf{x}^t | \mathbf{x}_0)} \left[\sum_{n=1}^L w(n, \mathbf{x}_{<n}^t) \log p_{\theta}(x_n | \mathbf{x}_{<n}^t) \right]. \quad (5)$$

This formulation generates dense supervision for the entire sequence in a single forward pass, theoretically preserving the $O(L)$ efficiency of standard ARMs without the computational overhead of block-wise vectorization.

However, strictly enforcing a causal constraint within a diffusion framework introduces a unique pathological state we term Information Collapse, which makes naive implementation unstable. In bidirectional architectures (e.g., BERT or MDLM), every token attends to the full global sequence. Even if a local region is heavily masked, the model can anchor its predictions on future tokens, maintaining a relatively uniform information density across positions. In contrast, under a causal mask, the visible context for a token x_n is strictly limited to its predecessors $\mathbf{x}_{<n}$. This creates a severe information asymmetry: early tokens with short histories are extremely vulnerable to corruption. For instance, if the first few tokens of a sequence are masked, predicting the subsequent token becomes mathematically equivalent to random guessing, as there is neither past history nor future context to rely on. Conversely, later tokens in long sequences often possess redundant history and remain predictable even under moderate noise.

Standard uniform diffusion strategies ignore this asymmetry, treating the blind guessing scenarios of early tokens equally with the well-supported predictions of later tokens. Forcing the model to minimize loss on these invalid contexts re-

sults in high-variance gradients and optimization instability. To make CARD effective, we must explicitly address this variable reliability of the causal context. We propose two complementary strategies: *Soft Tail Masking* (Section 3.2) to structurally guarantee that the historical context retains valid signals, and *Context-aware Reweighting* (Section 3.3) to adaptively down-weight predictions where the context remains too ambiguous.

3.2. Soft Tail Masking

Causal diffusion requires a noise strategy that respects the autoregressive nature of the model. Standard uniform masking is ill-suited here because it randomly corrupts tokens anywhere, including the sequence start ($n \ll L$). Since early tokens inherently possess little history, masking their few available context tokens effectively forces the model to predict from pure noise. To guarantee a valid historical context, a natural intuition is to concentrate all corruption at the sequence tail. This maximizes the clean prefix, ensuring stable supervision. However, strict tail masking completely removes the immediate neighbors of the corrupted tokens, ignoring the strong local dependencies required for language modeling (Khandelwal et al., 2018).

We propose *Soft Tail Masking* (Figure 3), a strategy designed to alleviate the issue by restricting masking to a dynamic tail window $[\max(0, L - \lambda t \cdot L), L]$. By maintaining a clean prefix while creating a mixed-state transition zone at the tail, we ensure the model accesses sufficient global history while retaining the local context needed for prediction. We prove in Appendix A (Proposition 2) that this preserves a higher lower bound on Mutual Information than uniform masking.

3.3. Context-aware Reweighting

Compared to standard ARMs, CARD predicts x_n from a stochastically corrupted prefix $\mathbf{x}_{< n}$. When the prefix is heavily masked, the conditional entropy $H(x_n | \mathbf{x}_{< n}^t)$ increases sharply. Forcing the model to produce confident predictions under such high uncertainty results in noisy gradients and optimization instability. The diffusion models typically employ a global weighting scheme (e.g., $1/t$ in MDLM) to balance contributions across noise levels at sequence level, grounded in the ELBO framework. However, global weighting is insufficient for causal models since the effective noise level varies *locally* at each token position n .

Considering the causal characteristics of CARD, we introduce a context-aware reweighting mechanism. Specifically, we propose to evaluate the ambiguity of the context $\mathbf{x}_{< n}^t$ along three dimensions: Quantity (total noise count), Distance (proximity of noise to target), and Density (consecutive corruption). These factors are synthesized into a unified local ambiguity score S_n^{local} , defined as the distance-

weighted sum of corruption costs in the history:

$$S_n^{local} = \sum_{i=1}^n C_i \cdot (1 - p)^{(n-i)}, \quad (6)$$

$$C_i = \mathbb{I}[x_i = [\text{MASK}]] \cdot (1 + \mathbb{I}[x_{i-1} = [\text{MASK}]]). \quad (7)$$

The formulation explicitly maps the three dimensions to mathematical components:

- **Noise Quantity:** The summation $\sum_{i=1}^n$ accumulates the corruption costs across the history. This term ensures that a higher total number of masked tokens leads to a larger cumulative score, naturally suppressing the weight for heavily corrupted contexts.
- **Noise Distance:** Following previous findings that the relevance of historical tokens decays exponentially with distance (Khandelwal et al., 2018; Lin & Tegmark, 2017), we introduce the decay factor $(1 - p)^{(n-i)}$, where p is a decay factor set to a constant 0.5. It ensures that noise in the immediate context will be penalized more heavily than noise in the distant past, as the immediate context is most critical for next-token prediction.
- **Noise Density:** The cost term C_i assigns a higher cost to consecutive masked tokens (e.g., spans), reflecting the difficulty of reconstructing regions where local dependencies are entirely severed.

Finally, the context-aware loss weight $w(n, \mathbf{x}_{< n}^t)$ is computed as:

$$w(n, \mathbf{x}_{< n}^t) = \frac{1}{\beta + S_n^{local}}, \quad (8)$$

where β is a smoothing constant (typically set to 1).

Our mechanism shifts the reweighting granularity from the sequence level (as in MDLM and BD3LM) to the token level. By down-weighting tokens in degraded contexts, the model focuses on regions with sufficient signal, leading to more efficient optimization (see Appendix A).

3.4. Confidence-Based Block Inference

We employ a confidence-based block sampling strategy to accelerate generation. Specifically, at each generation step, we initialize a candidate block of length K by appending mask tokens to the sequence tail, denoted as $\mathbf{x}^{(0)} = \{[\text{MASK}]_1, \dots, [\text{MASK}]_K\}$. We then perform iterative parallel denoising, where a token x_i at iteration j is updated only if its prediction probability exceeds a threshold

Table 1. LM Evaluation Harness results. All models are 1B parameters trained on 300B tokens.

| Model | ARC-Challenge 25-shot | ARC-Easy 25-shot | CommonsenseQA 7-shot | HellaSwag 3-shot | MMLU-redux 5-shot | PIQA 0-shot | SciQ 0-shot | Winogrande 5-shot | AVG |
|------------------------------|--------------------------|---------------------|-------------------------|---------------------|----------------------|----------------|----------------|----------------------|--------------|
| <i>Autoregressive Models</i> | | | | | | | | | |
| ARM | 34.04 | 64.65 | 52.74 | 61.26 | 25.45 | 75.95 | 81.10 | 55.96 | 56.39 |
| <i>Diffusion Models</i> | | | | | | | | | |
| BD3LM | 27.30 | 48.06 | 44.06 | 42.48 | 26.93 | 59.79 | 79.60 | 51.38 | 47.45 |
| MDLM | 29.44 | 49.16 | 36.45 | 48.32 | 26.36 | 59.63 | 76.60 | 54.46 | 47.55 |
| CARD (Ours) | 32.68 | 60.77 | 48.73 | 53.29 | 25.65 | 71.71 | 79.80 | 53.28 | 53.23 |

Table 2. PPL evaluation on various text domains. Lower is better.

| Model | AG News | arXiv | LAMBADA | LM1B | OpenWebText | PTB | PubMed | WikiText | AVG |
|------------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| <i>Autoregressive Models</i> | | | | | | | | | |
| ARM | 30.62 | 18.15 | 33.83 | 39.14 | 17.68 | 117.56 | 11.93 | 40.52 | 38.68 |
| <i>Diffusion Models</i> | | | | | | | | | |
| BD3LM | 41.18 | 44.60 | 39.17 | 40.04 | 40.97 | 118.30 | 34.66 | 39.28 | 49.78 |
| MDLM | 42.20 | 23.58 | 35.87 | 48.68 | 20.77 | 168.19 | 17.23 | 42.18 | 49.84 |
| CARD (Ours) | 27.67 | 20.34 | 30.36 | 29.61 | 17.59 | 97.74 | 13.20 | 38.67 | 34.40 |

τ (Wu et al., 2025; Cheng et al., 2025):

$$x_i^{(j+1)} = \begin{cases} x_i^{(j)} & \text{if } x_i^{(j)} \neq [\text{MASK}], \\ \arg \max_w p_\theta(w|\mathbf{x}^{(j)}) & \text{if } \max_w p_\theta(w|\mathbf{x}^{(j)}) > \tau, \\ [\text{MASK}] & \text{otherwise.} \end{cases} \quad (9)$$

To strictly bound latency, we impose a maximum step limit T_{max} . If the block is not fully denoised within T_{max} steps, all remaining masks are immediately decoded. Finally, the generated block is added to the KV cache. This approach allows the inference speed to be dynamically controlled by adjusting the block size K , threshold τ , and step limit T_{max} .

4. Experiments

To validate the effectiveness of the CARD framework, we benchmarked it against three architectures: ARM, MDLM, and BD3LM. All models were pre-trained on a 300B-token subset of FineWeb (Penedo et al., 2024) and aligned to a 1B-parameter scale. To ensure a fair comparison, the baselines utilized state-of-the-art optimizations: MDLM adopted variable-length packed QKV operators from Flash Attention (Dao, 2024), while BD3LM integrated `torch.compile` with Flex Attention (Dong et al., 2024). Detailed model structure configurations and training hyperparameters are provided in Appendix B.

4.1. Computational Efficiency

We first address the training cost bottleneck typical of diffusion models. Normalizing the training latency of ARM and CARD to a baseline of $1.0\times$, MDLM incurs a $1.5\times$ cost due to its bidirectional attention mechanism, while BD3LM

risers to roughly $3.0\times$ driven by input duplication constraints. In contrast, CARD eliminates these overheads, achieving superior performance while maintaining ARM-level training efficiency.

4.2. Performance Evaluation

Downstream Task Accuracy. We assessed disciplinary knowledge using ARC-Challenge & ARC-Easy (Clark et al., 2018), MMLU (Hendrycks et al., 2021), and SciQ (Welbl et al., 2017); commonsense reasoning via PIQA (Bisk et al., 2020), HellaSwag (Zellers et al., 2019), and CommonsenseQA (Talmor et al., 2019); and context disambiguation with Winogrande (Sakaguchi et al., 2020). As detailed in Table 1, a distinct performance hierarchy is evident. While the baseline methods MDLM and BD3LM plateau at an average of approximately 47.50%, CARD establishes significantly better results for non-autoregressive models with an average accuracy of 53.23%. The substantial 5.7% absolute improvement over prior diffusion baselines indicates that CARD effectively mitigates the performance degradation. Crucially, while a marginal gap to the autoregressive (ARM) upper bound remains, CARD significantly narrows this disparity, demonstrating that dense supervision can yield ARM-competitive performance without sacrificing the efficiency benefits of parallel decoding.

Language Modeling and Generalization. To evaluate intrinsic generative quality, we measured zero-shot perplexity across three distinct domains: general corpora using WikiText (Merity et al., 2017) and OpenWebText (Gokaslan et al., 2019); news and periodicals using AG News (Zhang et al., 2015), LM1B (Jozefowicz et al., 2016), and PTB (Marcus

Table 3. Perplexity (PPL) results on the LM1B dataset. Models are 110M parameters trained on 33B tokens using EMA.

| Model | ARM | BD3LM | MDLM | CARD |
|-------|-------|-------|-------|-------|
| PPL ↓ | 21.12 | 35.06 | 37.48 | 21.54 |

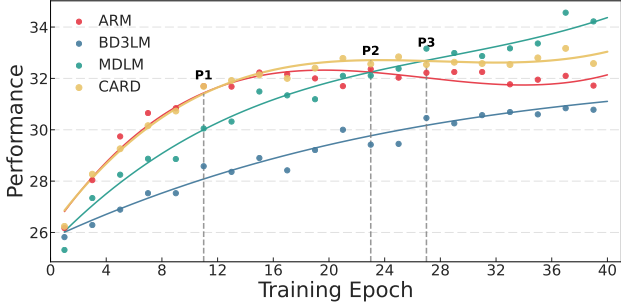


Figure 4. HellaSwag performance of four paradigms under repeated training on a FineWeb-Edu subset. The annotations mark specific crossover points in performance: P1 denotes the epoch where CARD surpasses ARM, P2 where MDLM overtakes ARM, and P3 where MDLM exceeds CARD.

et al., 1993); and specialized or long-context tasks using arXiv, PubMed (Cohan et al., 2018), and LAMBADA (Paperno et al., 2016). The results in Table 2 show that CARD consistently outperforms both diffusion baselines. More notably, CARD surpasses the ARM baseline on 6 out of 8 datasets. It achieves the best overall scores on general domains and the context-heavy LAMBADA benchmark, while remaining competitive on the specialized scientific vocabularies of arXiv and PubMed. We attribute the generalization advantage to the training objective. Standard ARMs rely on next-token prediction, which has been argued to be “myopic,” prioritizing local correlations and rote memorization (Nagarajan et al., 2025). Conversely, CARD’s denoising objective functions as a form of “teacherless training,” forcing the model to predict tokens from corrupted contexts. The mechanism incentivizes the model to capture global structural patterns and long-range dependencies rather than relying on local statistical shortcuts, resulting in superior generalization on unseen data compared to the strictly left-to-right of ARMs.

5. Analysis

5.1. Training Stability

A bottleneck in scaling discrete diffusion models is optimization instability. As derived in Appendix A (Proposition 3), BD3LM suffers from distributional discontinuities at block boundaries, while MDLM encounters high variance when predicting tokens from heavily masked contexts. To address these issues, practical implementations of these architectures often rely heavily on Exponential Moving Average

(EMA). Although rarely highlighted in their theoretical formulations, EMA is adopted by default in the official training repositories of both MDLM and BD3LM as a necessary stabilizer. In contrast, CARD is designed for inherent stability through its continuous causal loss landscape and context-aware reweighting (Proposition 1), which minimizes gradient variance by design. This theoretical guarantee reduces the dependency on aggressive parameter smoothing, ensuring that the optimization trajectory remains true to the underlying data distribution.

Empirical Validation with EMA. To investigate the training stability in detail, we conducted a controlled study on the LM1B dataset. We trained all models (110M parameters, 33B tokens) using the EMA configurations explicitly found in the baseline codebases. As shown in Table 3, even with EMA effectively buffering the gradient noise, MDLM and BD3LM yield perplexities of 37.48 and 35.06, respectively. In comparison, CARD achieves a significantly lower perplexity of 21.54 under identical conditions. The result demonstrates that while baselines require EMA to manage their structural instability, CARD utilizes it to further refine its density estimation, converging to a solution comparable to the ARM baseline.

5.2. Data Potential and Epoch Scaling

We define *Data Potential* as an architecture’s capacity to continuously extract signal from a fixed data distribution over repeated training iterations. Theoretically, based on the number of learnable conditional probability paths per datum (derived in Appendix C), we posit a hierarchy of $\text{MDLM} > \text{CARD} > \text{BD3LM} > \text{ARM}$, suggesting that ARMs saturate rapidly, whereas diffusion-based models sustain gains over longer horizons.

Empirical validation on 1B-parameter models trained on a 1B token subset of FineWeb-Edu (Penedo et al., 2024) over 40 epochs confirms the ranking (Figure 4). At the early training stage, CARD surpasses the ARM baseline at the inflection point P1 (\approx epoch 11) as the latter saturates. MDLM later overtakes ARM (P2) and eventually CARD (P3). Crucially, the interval preceding P3 identifies a functional “sweet spot”: CARD significantly outperforms ARM without requiring the extensive training horizon MDLM needs to realize its full potential.

The result has critical implications given the current scarcity of high-quality data, which necessitates training beyond Chinchilla-optimal ratios to minimize inference costs. While standard ARMs are ill-suited for the regime due to early saturation and MDLMs incur high initial compute costs, CARD effectively bridges the gap. It extends the performance boundary within practical computational budgets, offering a superior scaling solution when data quantity is the primary bottleneck.

Table 4. Ablation study on noise position and context-aware reweighting mechanisms.

| Setting | ARC-C 25-shot | ARC-E 25-shot | CSQA 7-shot | HellaS 3-shot | MMLU 5-shot | PIQA 0-shot | SciQ 0-shot | Wino 5-shot | AVG |
|-----------------------------------|------------------|------------------|----------------|------------------|----------------|----------------|----------------|----------------|--------------|
| CARD (Ours) | 32.68 | 60.77 | 48.73 | 53.29 | 25.45 | 71.71 | 79.80 | 53.28 | 53.21 |
| <i>Ablation on Noise Position</i> | | | | | | | | | |
| w/o Relaxed Window (Strict Tail) | 32.85 | 59.72 | 48.24 | 52.33 | 25.38 | 71.28 | 78.90 | 51.31 | 52.50 |
| w/o Tail Preference (Random) | 29.95 | 59.38 | 46.03 | 51.78 | 25.14 | 71.00 | 76.30 | 53.36 | 51.62 |
| <i>Ablation on Weighting</i> | | | | | | | | | |
| w/o Context-aware Weighting | 30.63 | 58.29 | 47.01 | 50.14 | 25.23 | 70.35 | 79.20 | 52.41 | 51.66 |

Table 5. Gen PPL results.

| Decoding Configuration | Throughput (tok/s) | AVG PPL ↓ |
|------------------------|--------------------|-----------|
| ARM (Baseline) | 3,771 (1.00×) | 11.19 |
| <i>CARD (Ours)</i> | | |
| Block=16, Steps=16 | 6,441 (1.71×) | 12.65 |
| Block=16, Steps=8 | 10,702 (2.84×) | 13.81 |
| Block=32, Steps=8 | 15,064 (4.01×) | 18.38 |

5.3. Ablation Study

In addition to the architectural comparison, we conducted ablation studies to validate the effectiveness of our proposed noise position preference (Section 3.2) and context-aware reweighting mechanisms (Section 3.3). The results are summarized in Table 4, leading to the following observations.

The noise distribution strategy plays a crucial role in unidirectional models. As shown in the results, applying noise to random positions (w/o Tail Preference) yields the lowest performance among the noise strategies. In a causal framework, tokens at the beginning of the sequence lack preceding context. If these tokens are masked randomly, the model cannot recover them effectively, leading to training inefficiencies. By concentrating noise at the tail, we observe a clear performance improvement. This suggests that a tail-biased noise strategy better aligns with the generative nature of language modeling, where history is used to predict the future. Furthermore, the results highlight the importance of the relaxed noise window. The “Strict Tail” setting, where the end of the sequence is a solid block of noise, underperforms compared to the full CARD implementation. A solid noise block creates an information void where the final tokens lack any immediate local context. By allowing a mix of clean and noisy tokens within the tail window (Relaxed Window), we enable the model to leverage local cues even during the denoising process.

Removing context-aware reweighting results in a noticeable drop in accuracy across most benchmarks. The dynamic weighting mechanism, rooted in the ELBO formulation, uses noise intensity to balance the training objective.

It naturally integrates the next-token prediction task with the diffusion objective by assigning appropriate importance to each token based on the clarity of its context. This ensures that the model focuses on learnable patterns rather than being overwhelmed by high-entropy predictions in heavily corrupted contexts.

5.4. Generation Perplexity Analysis

To further evaluate the generation quality, we conducted a generation perplexity (Gen PPL) analysis on Hellaswag prefixes using the model trained in our main experiment. For robust evaluation, we report the average PPL computed by four base models: Qwen3-8B (Qwen3 et al., 2025), SmolLM3-3B (Bakouch et al., 2025), gemma-3-27b (Gemma et al., 2025), and gpt2-large (Radford et al., 2019). All inference tests were performed with a batch size of 128. As shown in Table 5, our method demonstrates a promising trade-off between speed and quality. Specifically, we achieve a 1.62× speedup while maintaining a generation quality comparable to the ARM baseline. Furthermore, in a more aggressive setting, our method delivers over 4× inference acceleration with only a slight increase in PPL. These results strongly validate the potential of our method to serve as a new baseline for efficient generation. Additionally, we provide a detailed case study and discuss the potential failure modes of parallel generation in Appendix D.

6. Conclusion

We presented CARD, a unified framework that reconciles the training stability of autoregressive models with the parallel inference capabilities of diffusion. By reformulating discrete diffusion within a strict causal constraint, CARD eliminates the computational overhead of block-based architectures. Empirically, CARD not only matches the generation quality of standard ARMs but also speed up to 1.7× through dynamic parallel decoding. Crucially, our analysis of data potential reveals that CARD avoids early saturation in multi-epoch regimes, positioning it as a highly data-efficient backbone for next-generation LLMs.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning, particularly by improving the training and inference efficiency of Large Language Models. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Arriola, M., Sahoo, S. S., Gokaslan, A., Yang, Z., Qi, Z., Han, J., Chiu, J. T., and Kuleshov, V. Block diffusion: Interpolating between autoregressive and diffusion language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=tyEyYT267x>.
- Austin, J., Johnson, D. D., Ho, J., Tarlow, D., and van den Berg, R. Structured denoising diffusion models in discrete state-spaces. In *Proceedings of the 35th International Conference on Neural Information Processing Systems*, NIPS '21, Red Hook, NY, USA, 2021. Curran Associates Inc. ISBN 9781713845393.
- Bakouch, E. et al. SmolLM3: smol, multilingual, long-context reasoner. <https://huggingface.co/blog/smollm3>, 2025.
- Bie, T., Cao, M., Chen, K., Du, L., Gong, M., Gong, Z., Gu, Y., Hu, J., Huang, Z., Lan, Z., Li, C., Li, C., Li, J., Li, Z., Liu, H., Liu, L., Lu, G., Lu, X., Ma, Y., Tan, J., Wei, L., Wen, J.-R., Xing, Y., Zhang, X., Zhao, J., Zheng, D., Zhou, J., Zhou, J., Zhou, Z., Zhu, L., and Zhuang, Y. Llada2.0: Scaling up diffusion language models to 100b, 2025. URL <https://arxiv.org/abs/2512.15745>.
- Bisk, Y., Zellers, R., Bras, R. L., Gao, J., and Choi, Y. PIQA: reasoning about physical commonsense in natural language. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pp. 7432–7439. AAAI Press, 2020. doi: 10.1609/AAAI.V34I05.6239. URL <https://doi.org/10.1609/aaai.v34i05.6239>.
- Cheng, S., Bian, Y., Liu, D., Zhang, L., Yao, Q., Tian, Z., Wang, W., Guo, Q., Chen, K., Qi, B., and Zhou, B. Sdar: A synergistic diffusion-autoregression paradigm for scalable sequence generation, 2025. URL <https://arxiv.org/abs/2510.06303>.
- Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., and Tafjord, O. Think you have solved question answering? try arc, the ai2 reasoning challenge, 2018. URL <https://arxiv.org/abs/1803.05457>.
- Cohan, A., Dernoncourt, F., Kim, D. S., Bui, T., Kim, S., Chang, W., and Goharian, N. A discourse-aware attention model for abstractive summarization of long documents. In Walker, M., Ji, H., and Stent, A. (eds.), *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pp. 615–621, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2097. URL <https://aclanthology.org/N18-2097/>.
- Dao, T. Flashattention-2: Faster attention with better parallelism and work partitioning. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=mZn2Xyh9Ec>.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. BERT: Pre-training of deep bidirectional transformers for language understanding. In Burstein, J., Doran, C., and Solorio, T. (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423/>.
- Dong, J., Feng, B., Guessous, D., Liang, Y., and He, H. Flex attention: A programming model for generating optimized attention kernels, 2024. URL <https://arxiv.org/abs/2412.05496>.
- Gemma et al. Gemma 3 technical report, 2025. URL <https://arxiv.org/abs/2503.19786>.
- Gokaslan, A., Cohen, V., Pavlick, E., and Tellex, S. Openwebtext corpus. <http://Skylion007.github.io/OpenWebTextCorpus>, 2019.
- Gong, S., Li, M., Feng, J., Wu, Z., and Kong, L. Dif-fuseq: Sequence to sequence text generation with diffusion models. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=jQj-_rLVXsj.
- Han, K., Shan, N., Zhao, Z., Hu, Z., Dong, X., Ye, J., Pan, L., Wu, F., and Kuang, K. C²dml: Causal concept-guided diffusion large language models, 2025. URL <https://arxiv.org/abs/2511.22146>.

- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=d7KBjmI3GmQ>.
- Jozefowicz, R., Vinyals, O., Schuster, M., Shazeer, N., and Wu, Y. Exploring the limits of language modeling, 2016. URL <https://arxiv.org/abs/1602.02410>.
- Khandelwal, U., He, H., Qi, P., and Jurafsky, D. Sharp nearby, fuzzy far away: How neural language models use context. In Gurevych, I. and Miyao, Y. (eds.), *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 284–294, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1027. URL <https://aclanthology.org/P18-1027/>.
- Kim, J., Shah, K., Kontonis, V., Kakade, S. M., and Chen, S. Train for the worst, plan for the best: Understanding token ordering in masked diffusions. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=DjJmre5IkP>.
- Li, X., Thickstun, J., Gulrajani, I., Liang, P. S., and Hashimoto, T. B. Diffusion-lm improves controllable text generation. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 4328–4343. Curran Associates, Inc., 2022.
- Lin, H. W. and Tegmark, M. Critical behavior in physics and probabilistic formal languages. *Entropy*, 19(7), 2017. ISSN 1099-4300. doi: 10.3390/e19070299. URL <http://www.mdpi.com/1099-4300/19/7/299>.
- Liu, A., He, M., Zeng, S., Zhang, S., Zhang, L., Wu, C., Jia, W., Liu, Y., Zhou, X., and Zhou, J. Wedlm: Reconciling diffusion language models with standard causal attention for fast inference, 2025. URL <https://arxiv.org/abs/2512.22737>.
- Lou, A., Meng, C., and Ermon, S. Discrete diffusion modeling by estimating the ratios of the data distribution. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=CNicRIVIPA>.
- Marcus, M. P., Santorini, B., and Marcinkiewicz, M. A. Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2):313–330, 1993. URL <https://aclanthology.org/J93-2004/>.
- Merity, S., Xiong, C., Bradbury, J., and Socher, R. Pointer sentinel mixture models. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=Byj72udxe>.
- Nagarajan, V., Wu, C. H., Ding, C., and Raghunathan, A. Roll the dice & look before you leap: Going beyond the creative limits of next-token prediction. In *Forty-second International Conference on Machine Learning, ICML 2025, Vancouver, BC, Canada, July 13-19, 2025*. OpenReview.net, 2025. URL <https://openreview.net/forum?id=Hi0SyHMmkd>.
- Ni, J., Liu, Q., Dou, L., Du, C., Wang, Z., Yan, H., Pang, T., and Shieh, M. Q. Diffusion language models are super data learners, 2025. URL <https://arxiv.org/abs/2511.03276>.
- Nie, S., Zhu, F., Du, C., Pang, T., Liu, Q., Zeng, G., Lin, M., and Li, C. Scaling up masked diffusion models on text. In *The Thirteenth International Conference on Learning Representations*, 2025a. URL <https://openreview.net/forum?id=WNvwwK0tut>.
- Nie, S., Zhu, F., You, Z., Zhang, X., Ou, J., Hu, J., ZHOU, J., Lin, Y., Wen, J.-R., and Li, C. Large language diffusion models. In *ICLR 2025 Workshop on Deep Generative Model in Machine Learning: Theory, Principle and Efficiency*, 2025b. URL <https://openreview.net/forum?id=wzl6ltIUj6>.
- Ou, J., Nie, S., Xue, K., Zhu, F., Sun, J., Li, Z., and Li, C. Your absorbing discrete diffusion secretly models the conditional distributions of clean data. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=sMyXP8Tanm>.
- Paperno, D., Kruszewski, G., Lazaridou, A., Pham, Q., Bernardi, R., Pezzelle, S., Baroni, M., Boleda, G., and Fernández, R. The lambada dataset: Word prediction requiring a broad discourse context. volume 3, pp. 10. Association for Computational Linguistics (ACL), 2016. ISBN 9781510827585. doi: 10.18653/v1/p16-1144.
- Penedo, G., Kydlíček, H., allal, L. B., Lozhkov, A., Mitchell, M., Raffel, C., Von Werra, L., and Wolf, T. The fineweb datasets: Decanting the web for the finest text data at scale. In Globerson, A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J., and Zhang, C. (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 30811–30849. Curran Associates, Inc., 2024. doi: 10.52202/079017-0970.
- Qwen3 et al. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.

- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. Language Models are Unsupervised Multitask Learners. 2019. URL <https://openai.com/blog/better-language-models/>.
- Sahoo, S. S., Arriola, M., Schiff, Y., Gokaslan, A., Marroquin, E., Chiu, J. T., Rush, A., and Kuleshov, V. Simple and effective masked diffusion language models. In Globerson, A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J., and Zhang, C. (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 130136–130184. Curran Associates, Inc., 2024. doi: 10.52202/079017-4135.
- Sakaguchi, K., Bras, R. L., Bhagavatula, C., and Choi, Y. Winogrande: An adversarial winograd schema challenge at scale. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pp. 8732–8740. AAAI Press, 2020. doi: 10.1609/AAAI.V34I05.6399. URL <https://doi.org/10.1609/aaai.v34i05.6399>.
- Shi, J., Han, K., Wang, Z., Doucet, A., and Titsias, M. Simplified and generalized masked diffusion for discrete data. In Globerson, A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J., and Zhang, C. (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 103131–103167. Curran Associates, Inc., 2024. doi: 10.52202/079017-3277.
- Talmor, A., Herzig, J., Lourie, N., and Berant, J. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In Burstein, J., Doran, C., and Solorio, T. (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4149–4158, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1421. URL <https://aclanthology.org/N19-1421/>.
- Wang, G., Schiff, Y., Sahoo, S. S., and Kuleshov, V. Re-masking discrete diffusion models with inference-time scaling. In *ICLR 2025 Workshop on Deep Generative Model in Machine Learning: Theory, Principle and Efficacy*, 2025. URL <https://openreview.net/forum?id=xNwZ8kDC7T>.
- Welbl, J., Liu, N. F., and Gardner, M. Crowdsourcing multiple choice science questions. In Derczynski, L., Xu, W., Ritter, A., and Baldwin, T. (eds.), *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pp. 94–106, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/W17-4413. URL <https://aclanthology.org/W17-4413/>.
- Wu, C., Zhang, H., Xue, S., Liu, Z., Diao, S., Zhu, L., Luo, P., Han, S., and Xie, E. Fast-dllm: Training-free acceleration of diffusion llm by enabling kv cache and parallel decoding, 2025. URL <https://arxiv.org/abs/2505.22618>.
- Ye, J., Xie, Z., Zheng, L., Gao, J., Wu, Z., Jiang, X., Li, Z., and Kong, L. Dream 7b: Diffusion large language models, 2025. URL <https://arxiv.org/abs/2508.15487>.
- Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi, Y. HellaSwag: Can a machine really finish your sentence? In Korhonen, A., Traum, D., and Màrquez, L. (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4791–4800, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1472. URL <https://aclanthology.org/P19-1472/>.
- Zhang, X., Zhao, J., and LeCun, Y. Character-level convolutional networks for text classification. In Cortes, C., Lawrence, N., Lee, D., Sugiyama, M., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. URL https://proceedings.neurips.cc/paper_files/paper/2015/file/250cf8b51c773f3f8dc8b4be867a9a02-Paper.pdf.

A. Mathematical Foundations of CARD

In this section, we provide a formal analysis of the optimization dynamics and information-theoretic properties of the Causal Autoregressive Diffusion (CARD) framework. We contrast CARD with Masked Discrete Diffusion Models (MDLM) and Block-wise Discrete Diffusion Models (BD3LM).

A.1. Notation and Preliminaries

Let $\mathbf{x} = (x_1, \dots, x_L)$ be a sequence of length L from a discrete vocabulary \mathcal{V} . Let $\mathcal{M} \subset \{1, \dots, L\}$ denote the set of indices masked at time $t \in [0, 1]$. For any position n , we define the *causal context* $\mathcal{C}_n = \{x_i \mid i < n, i \notin \mathcal{M}\}$. The training objective is to minimize the expected negative log-likelihood:

$$\mathcal{L}(\theta) = \mathbb{E}_{t, \mathcal{M}} \left[\sum_{n=1}^L w(n, \mathcal{C}_n) \cdot \ell_n(\theta; \mathcal{C}_n) \right] \quad (10)$$

where $\ell_n(\theta; \mathcal{C}_n) = -\log p_\theta(x_n \mid \mathcal{C}_n)$ and $w(n, \mathcal{C}_n)$ is the weight assigned to the prediction at position n .

Definition A.1 (Local Ambiguity Score). The Local Ambiguity Score S_n^{local} is defined as a weighted sum of corruption costs within the causal window:

$$S_n^{local}(\mathcal{C}_n) = \sum_{i=1}^{n-1} C_i \cdot (1-p)^{n-i} \quad (11)$$

where $C_i = \mathbb{I}[i \in \mathcal{M}] \cdot (1 + \mathbb{I}[i-1 \in \mathcal{M}])$ represents the cost of masking, and $p \in (0, 1)$ is a decay factor.

A.2. Proposition 1: Gradient Variance Stabilization

Proposition A.2. The CARD weighting scheme $w(n, \mathcal{C}_n) = (\beta + S_n^{local})^{-1}$ minimizes the variance of the stochastic gradient estimator by performing an instance-level inverse-variance weighting.

Proof. Consider the variance of the stochastic gradient $\mathbf{g}_n = \nabla_\theta \ell_n$. In the discrete diffusion setting, as the context \mathcal{C}_n becomes increasingly corrupted (high S_n^{local}), the conditional distribution $p_\theta(x_n \mid \mathcal{C}_n)$ approaches the uninformative marginal distribution $p(x_n)$. In this regime, the Fisher Information $\mathcal{I}(\theta)_n = \mathbb{E}[\nabla_\theta \ell_n \nabla_\theta \ell_n^\top]$ is dominated by the noise of the sampling process rather than the underlying structural signal of the language.

Let $\sigma_n^2(\mathcal{C}_n) = \|\nabla_\theta \ell_n(\mathcal{C}_n)\|^2$ be the squared norm of the gradient. Given the power-law decay of mutual information in sequences, we posit that σ_n^2 is monotonically bounded by the ambiguity score: $\sigma_n^2 \leq \alpha S_n^{local} + \epsilon$. The variance of the weighted estimator is:

$$\text{Var}[w \cdot \mathbf{g}_n] = \mathbb{E}[w^2 \|\mathbf{g}_n\|^2] - \|\mathbb{E}[w \mathbf{g}_n]\|^2 \leq \frac{\alpha S_n^{local} + \epsilon}{(\beta + S_n^{local})^2} \quad (12)$$

As $S_n^{local} \rightarrow \infty$, the weighted gradient norm $\|w \mathbf{g}_n\| \rightarrow 0$. This ensures that uninformative, high-entropy contexts do not contribute disproportionately to the parameter updates, satisfying the conditions for stable convergence in the absence of aggressive Exponential Moving Average (EMA). \square

A.3. Proposition 2: Signal Retention via Causal MI Maximization

Proposition A.3. For a fixed noise budget t , the Soft Tail Masking strategy preserves a strictly higher lower bound on the cumulative Mutual Information (MI) compared to Uniform Masking.

Proof. Let $I(x_n; x_i)$ be the MI between tokens. In natural language, $I(x_n; x_i) \approx f(|n-i|)$, where f is a monotonically decreasing function. The total information available to the model is $\mathcal{I}_{total} = \sum_{n=1}^L \sum_{i < n, i \notin \mathcal{M}} I(x_n; x_i)$.

1. **Uniform Masking:** For MDLM, each $i \in \mathcal{M}$ with probability t . The expected MI at position n is $(1-t) \sum_{i < n} I(x_n; x_i)$.
2. **Soft Tail Masking:** CARD restricts masks to the tail window. For $n < L(1-\lambda t)$, the probability $P(i \in \mathcal{M} \mid i < n) = 0$.

Since $I(x_n; x_i)$ is maximal when $n-i$ is small, the Soft Tail strategy ensures that for a significant portion of the sequence (the ‘‘Head’’), the model observes the full causal signal. Because $\sum_{n=1}^L I(x_n; \mathcal{C}_n^{CARD})$ prioritizes preserving low-distance dependencies which contain the highest MI, it follows that $\mathcal{I}_{total}^{CARD} > \mathcal{I}_{total}^{MDLM}$. \square

A.4. Proposition 3: Landscape Continuity and Block Discontinuity

Proposition A.4. *CARD eliminates the $O(1)$ distributional shift discontinuities present in block-wise diffusion architectures (BD3LM).*

Proof. Let μ_n be the distribution of the context \mathcal{C}_n . We evaluate the continuity of the loss landscape by the Total Variation (TV) distance between adjacent context distributions $d_{TV}(\mu_n, \mu_{n+1})$.

In **BD3LM**, sequences are partitioned into blocks $\{B_k\}$. At a boundary index j where $x_j \in B_k$ and $x_{j+1} \in B_{k+1}$, the context shifts from a deterministic clean history (from previous blocks) to a stochastic noisy context (within the current block). This implies:

$$\lim_{L \rightarrow \infty} d_{TV}(\mu_j, \mu_{j+1}) = \|p(x_{clean}) - p(x_{noisy})\|_{TV} \approx O(1) \quad (13)$$

This jump results in a non-Lipschitz gradient spike at every block boundary.

In **CARD**, the transition probability $P(x_n = [\text{MASK}])$ is defined by a continuous noise schedule $\sigma(n, t)$ over the sequence index. For a linear schedule, the change in masking probability between n and $n+1$ is $O(1/L)$. Thus, $d_{TV}(\mu_n, \mu_{n+1}) \leq \frac{K}{L}$, ensuring that the expected loss and its gradients are Lipschitz continuous with respect to the sequence index. \square

B. Experimental Setups

We detail the model architecture and training hyperparameters used in our experiments, with the full configuration summarized in Table 6.

Model Architecture Our model is built upon a bidirectional Transformer encoder architecture, incorporating Flash Attention 2 for computational efficiency. It consists of 33 Transformer layers with a hidden dimension of 1536 and an intermediate FFN dimension of 4096, utilizing the SiLU activation function. The model supports a maximum position embedding length of 8192 tokens.

Training Configuration Training is performed using the AdamW optimizer with `bfloat16` mixed precision. We employ a constant learning rate schedule with a 2,500-step warmup, peaking at 3×10^{-4} . For the diffusion process, the masking probability is linearly annealed from 1.0 to 0.

Table 6. Experimental Setup: Model Architecture and Training Hyperparameters

| Model Architecture | | Training Hyperparameters | |
|---------------------|--------------|--------------------------|--------------------------------|
| Parameter | Value | Hyperparameter | Value |
| Number of Layers | 33 | Optimizer | AdamW |
| Hidden Size | 1536 | Peak Learning Rate | 3×10^{-4} |
| Intermediate Size | 4096 | LR Scheduler | Cosine w/ Warmup |
| Attention Heads | 24 | Warmup Steps | 2,500 |
| Vocab Size | 50,368 | Max Training Steps | 1,000,000 |
| Activation Function | SiLU | Sequence Length | 128 |
| Max Pos Embeddings | 8,192 | Precision | BF16 |
| Attn Implementation | Flash Attn 2 | Noise Schedule | Linear (1.0 \rightarrow 0.0) |

C. Complexity Analysis of Learnable Conditional Probabilities

In this section, we quantify the number of structural conditional probabilities that different generative models can learn. We define L as the sequence length. We analyze the theoretical upper bound of dependency patterns based on the attention mechanism and the masking strategy employed by each model.

C.1. Autoregressive Models (ARM)

Standard Autoregressive Models rely on the probability chain rule. The generation of a token at position t depends strictly on the fixed sequence of preceding tokens x_1, \dots, x_{t-1} . Since the context for every position is deterministic and unique (the prefix), the model does not learn from varying subsets of the context. Therefore, the total number of learnable conditional probabilities is linear with respect to the sequence length:

$$N_{\text{ARM}} = L \quad (14)$$

C.2. Causal Autoregressive Diffusion (CARD)

CARD combines unidirectional attention with a discrete diffusion process. Although the attention mechanism restricts information flow from left to right, the noise injection process introduces combinatorial diversity. For a token at position t , the context consists of tokens x_1 to x_{t-1} . In the diffusion training process, each of these context tokens can exist in two states: masked or unmasked.

This results in a geometric series where the first token has 1 possible context state, the second has 2, and the last has 2^{L-1} . The total number of combinations is the sum of this series:

$$N_{\text{CARD}} = \sum_{t=0}^{L-1} 2^t = 2^L - 1 \quad (15)$$

C.3. Masked Discrete Language Models (MDLM)

MDLM represents the standard bidirectional discrete diffusion approach. The model utilizes bidirectional attention, allowing any token to attend to any other token in the sequence. During training, a random proportion of tokens are masked.

For any given target position i , the context is a subset of the remaining $L - 1$ tokens. Since each of the other tokens can be either masked or unmasked, there are 2^{L-1} possible context configurations for a single position. Since all L positions serve as prediction targets, the total number of learnable probabilities is:

$$N_{\text{MDLM}} = L \times 2^{L-1} \quad (16)$$

C.4. Blockwise Diffusion (BD3LM)

BD3LM employs a hybrid architecture. It divides the sequence of length L into N blocks, where each block has a size of K (such that $L = N \times K$). The model applies unidirectional causal attention between blocks but maintains bidirectional attention within each block.

Since the inter-block connection is causal, previous blocks act as a fixed context and do not contribute to combinatorial explosion. However, within each block of size K , the model behaves like a bidirectional diffusion model. The number of combinations per block is $K \times 2^{K-1}$. Summing this over all N blocks yields:

$$N_{\text{BD3LM}} = \frac{L}{K} \times (K \times 2^{K-1}) = L \times 2^{K-1} \quad (17)$$

C.5. Summary

Table 7 summarizes the number of learnable conditional probabilities for each model. This comparison highlights that while diffusion-based models offer exponentially larger state spaces than ARM, Blockwise Diffusion (BD3LM) effectively bridges the gap by controlling the exponent through the block size K .

D. Case Study: Impact of Acceleration Ratios

Table 8 presents a qualitative comparison between the ARM baseline and our CARD method. At a moderate acceleration ratio of $1.7\times$, CARD maintains generation quality comparable to the baseline, producing coherent and contextually appropriate text. However, aggressively increasing the speedup to $4\times$ by restricting the step budget leads to noticeable degradation. Instead of syntactic errors, this degradation primarily manifests as **logical repetition and text looping** (e.g., repeating

Table 7. Comparison of Learnable Conditional Probabilities

| Model Type | Attention Type | Complexity |
|------------|----------------|--------------------|
| ARM | Unidirectional | L |
| CARD | Unidirectional | $2^L - 1$ |
| BD3LM | Hybrid | $L \times 2^{K-1}$ |
| MDLM | Bidirectional | $L \times 2^{L-1}$ |

similar sentence structures or phrases). This phenomenon stems from the hard step limit: the model is compelled to complete the text block via non-autoregressive generation at the final step. Lacking sufficient autoregressive guidance, the model tends to collapse into high-probability repetitive patterns rather than developing diverse narrative progressions. While we anticipate that stronger base models will mitigate this sensitivity, we currently recommend adhering to the standard configuration to strike the optimal balance between speed and quality.

Table 8. Comparison of generation quality under different acceleration settings. **Case 1** demonstrates how aggressive speedup ($4\times$) leads to repetitive sentence structures. **Case 2** further illustrates it, where CARD ($4\times$) falls into a degenerative loop (repeating “applying the gel... bottle is shown”), whereas the baseline and moderate settings maintain narrative flow.

| Case Study 1 | |
|--------------------------------------|---|
| Prompt | <i>A man is sitting on a roof. he...</i> |
| ARM | sees a beautiful young lady on a rooftop. The man: a woman? The woman: no. The man: Oh! Well, I’m sure she’s beautiful. The woman: No, she’s not... |
| CARD (1.7\times) | ’s living a peaceful life, working hard to support his family and doing what he loves. But there are some days that he can’t do it anymore. He’s getting older, he’s getting tired... |
| CARD (4\times) | is looking at the sky. he is thinking about his life. he is thinking about his past. he is thinking about his future... |
| Case Study 2 | |
| Prompt | <i>A bottle of deep cleansing gel is shown. a woman</i> |
| ARM | ’s face is washed with a mild cleanser and she applies face lotion with aloe vera. a woman’s face is washed with a mild cleanser and she applies face lotion with aloe vera. |
| CARD (1.7\times) | is shown with legs raised. A bottle of liquid hand soap is shown. a woman is shown leaving the shop and a young man is seen patting her down. |
| CARD (4\times) | is shown applying the gel to her face. the bottle of the gel is shown. the woman is shown applying the gel to her face. the bottle of the gel is shown. A bottle of deep cleansing gel is shown. a woman is shown applying the gel to her face. the bottle of the gel is shown... |