

DON'T SETTLE TOO EARLY: SELF-REFLECTIVE REMASKING FOR DIFFUSION LANGUAGE MODELS

Zemin Huang^{1,2} Yuhang Wang² Zhiyang Chen^{2,4*} Guo-Jun Qi^{2,3,4,*}

¹ Zhejiang University

² MAPLE Lab, Westlake University

³ Matterwave Intelligence

⁴ Institute of Advanced Technology, Westlake Institute for Advanced Study
 {huangzemin, wangyuhang, chenzyang}@westlake.edu.cn,
 guojunq@gmail.com

<https://github.com/maple-research-lab/RemeDi>

ABSTRACT

Mask-based Diffusion Language Models (DLMs) struggle to revise incorrect tokens: once a token is generated, it typically remains fixed. The key challenge is to identify potential errors in the inputs. In this paper, we propose *Remasking-enabled Diffusion Language Model (RemeDi)*, a mask-based DLM that introduces *remasking* as another fundamental mechanism, enabling more flexible text refinement in diffusion-based text generation. To achieve this, RemeDi jointly predicts token distributions and per-token confidence scores at each step. The confidence scores determine which tokens to be unmasked after the current step, allowing the model to identify tokens with low quality and remask them. These remasked tokens can be resampled with richer context in subsequent steps. We design a remask-aware pipeline to train this ability, including supervised fine-tuning which teaches the model to detect and remask incorrect tokens in addition to predict mask tokens, and reinforcement learning which optimizes full generation trajectories toward higher rewards. Experiments show that RemeDi achieves the state-of-the-art results among open-source DLMs on multiple datasets.

1 INTRODUCTION

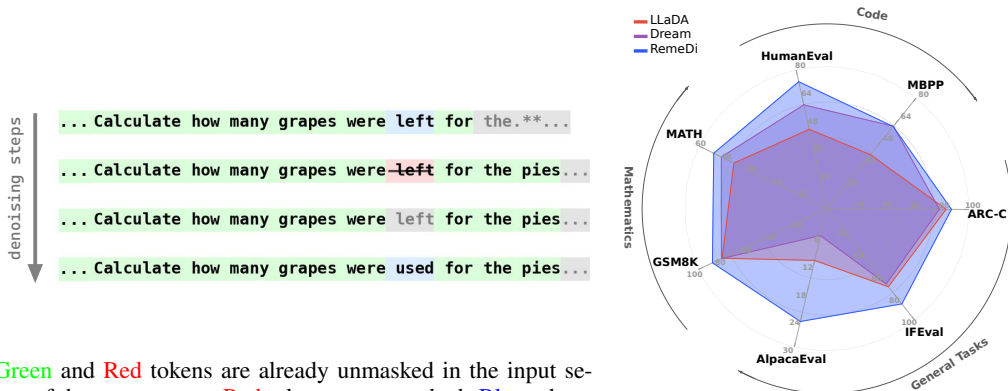
Diffusion Language Models (DLMs) have recently emerged as a promising alternative to autoregressive language models (Nie et al., 2025; Ye et al., 2025; Lou et al., 2024; Arriola et al., 2025). A DLM defines a forward process that gradually corrupts text into a noise prior, and learns a reverse process to recover clean text (Campbell et al., 2022; Lou et al., 2024). Unlike autoregressive models, DLMs do not commit to a fixed left-to-right order, offering greater flexibility in generation and an inherent ability to predict multiple tokens in parallel.

A dominant variant is the *mask-based* DLM (Nie et al., 2025; Ye et al., 2025), where the noise is represented by a special mask token. Under this formulation, the model learns to recover masked tokens during training, while assuming that once tokens are unmasked, they are supposed to be correct without having to clean them later. This assumption is problematic: the model may generate wrong tokens, which should be revealed and corrected in later steps when more contexts are available. However, most existing DLMs (Nie et al., 2025; Ye et al., 2025) keep already unmasked tokens fixed, preventing them from being revised by self-reflecting on errors.

To address this, several works have explored methods to revise generated tokens. von Rütte et al. (2025) defines a new noise schedule by interpolating between masking and uniform noise, enabling revision of wrong tokens on small-scale models. Wang et al. (2025a) applies predictor-corrector samplers by stochastically remasking a subset of tokens only at inference time, where the remasking is performed randomly without training the model how to find and remask incorrect tokens. For large-scale DLMs, Seed Diffusion (Song et al., 2025) allows all tokens to be resampled at every step.

*Corresponding authors

This project was initiated and supported by MAPLE Lab at Westlake University.



(a) Green and Red tokens are already unmasked in the input sequence of the current step. Red tokens are remasked. Blue tokens are unmasked in the outputs. Gray tokens keep masked in the outputs, and we display the token with the highest probability at these positions.

(b) Radar plot comparing the performance of RemeDi with other DLMs across various evaluation benchmarks.

Figure 1: (a) Illustration of remasking for quality improvement: RemeDi initially predicts the token “left”, but later identifies a semantic mismatch in the phrase “left for the pies”. The model then remasks this token and corrects it to the more appropriate “used”. (b) RemeDi outperforms existing DLMs in various tasks, including math, code and general benchmark.

However, it lacks a mechanism to ensure that the number of mask tokens decreases monotonically — it is a key feature for diffusion models to ensure decreasing noise levels over steps (Guo et al., 2025), so that the mask tokens will eventually vanish at the final step to complete the generation.

In this paper, we propose a self-reflective remasking approach to train DLMs. As illustrated in Fig. 1a, it aims to train DLMs with the ability of finding wrong tokens and turning them back to mask ones so that they can be resampled with richer context in later steps. Based on this, we introduce *Remasking-enabled Diffusion Language Model (RemeDi)*, a mask-based DLM that incorporates self-reflective remasking to revise already generated but incorrect tokens. RemeDi jointly predicts token distributions and per-token confidence scores. At each diffusion step, high-confidence tokens are unmasked while low-confidence ones are (re-)masked, regardless of whether they have been previously unmasked.

The key challenge is to *train* the model how to remask incorrect tokens in a self-reflective manner. To this end, we design a remask-aware training pipeline in two stages: 1) **Remask SFT**, where the model learns to identify and remask incorrect tokens, while predicting masked tokens. We construct an input sequence for Remask SFT by randomly masking its tokens or replacing them with random alternatives to simulate the noise. The noise schedule deciding how many tokens are masked or randomly replaced is designed to follow the criterion that the noise level should monotonically decrease over steps. The model is then trained to remask and revise incorrect tokens over noisy input sequences. 2) **Remask RL**, where the model is further fine-tuned with outcome-based reinforcement learning. It seeks to optimize the entire generation trajectories toward final outputs with higher rewards by considering how to remask and predict tokens in each step.

As shown in Fig. 1b, RemeDi achieves the state-of-the-art performance among open-source DLMs, achieving competitive results on various benchmark datasets, including math problems (89.1% on GSM8K (Cobbe et al., 2021), 52.9% on MATH (Hendrycks et al., 2021)), code generation (73.2% on HumanEval (Chen et al., 2021), 59.4% on MBPP (Austin et al., 2021)), and general tasks (24.5% on AlpacaEval (Dubois et al., 2024), 85.4% on IFEval (Zhou et al., 2023), and 87.7% on ARC-C (Clark et al., 2018)).

2 RELATED WORK

2.1 MASK-BASED DIFFUSION LANGUAGE MODELS

Diffusion language models (DLMs) have emerged as promising alternatives to auto-regressive (AR) models for text generation. Among them, mask-based DLMs (Nie et al., 2025; Ye et al., 2025; Zheng et al., 2023; Ou et al., 2024) dominate, generating text by progressively denoising mask tokens. Recent studies (Arriola et al., 2025; Fathi et al., 2025; Huang & Tang, 2025; Sahoo et al., 2025; Wang et al., 2025b; Gat et al., 2025) have increasingly explored the fusion of AR and diffusion models, often through an iterative block-wise decoding strategy: inference proceeds by iteratively appending a block of mask tokens to the input sequence and denoising it, repeating until the EOS token is generated. This paradigm inherits the strengths — flexible generation order and parallel decoding from DLMs, and cache efficiency from AR — yielding faster inference without sacrificing quality. In our work, we adapt LLaDA-8B-Instruct (Nie et al., 2025) to variable-length block-by-block generation, serving as the backbone for our remasking mechanism.

2.2 REVISING ERRORS IN DIFFUSION LANGUAGE MODELS

A key limitation of mask-based DLMs is their inability to revise tokens once unmasked, even if they are incorrect. Existing efforts to address this fall into two categories. The first category (Campbell et al., 2022; Wang et al., 2025a) applies predictor-corrector samplers without training, for example by stochastically remasking a subset of tokens during inference. These methods lack a mechanism to identify which tokens are actually wrong. As a result, they have to rely on many extra sampling steps to take effect, which are inefficient and hard to optimize. The second category modifies the diffusion process to enable revision during the reverse diffusion process, e.g., combining mask diffusion process with either the uniform diffusion process (von Rütte et al., 2025) or edit-based diffusion process (Havasi et al., 2025; Song et al., 2025).

In short, none of these approaches provides a principled way to detect and selectively correct errors during generation. In contrast, RemeDi fulfills self-reflection by identifying and remasking error-prone tokens through a two-stage learning pipeline, and jointly training the model to resample the remasked tokens in later steps.

3 METHODS

3.1 PRELIMINARIES: MASK-BASED DIFFUSION LANGUAGE MODELS

Diffusion Language Models (DLMs) aim to model text generation by approximating the probability distribution p_{data} over a finite vocabulary $\mathcal{V} = \{1, 2, \dots, V\}$. They define a discrete diffusion process in which the unknown data distribution p_{data} at $t = 0$ gradually evolves into a simple prior distribution p_{prior} at $t = T$ (Lou et al., 2024). At intermediate times t , we denote the distribution as p_t . Formally, this diffusion process can be described by a linear ODE involving a diffusion matrix Q_t :

$$\frac{dp_t}{dt} = Q_t p_t, \quad p_0 = p_{\text{data}}, \quad p_T = p_{\text{prior}}. \quad (1)$$

While t can be defined continuously as in Eq. (1), in practice we work with discrete timesteps $t_{0:N}$.

In this paper, we focus on mask-based DLMs, where p_{prior} is a distribution that puts all its mass on the mask state, denoted as [M]. Given a clean sequence $x_0 \sim p_{\text{data}}$, a corrupted sequence x_t is obtained by randomly replacing part of the tokens with the mask token [M]. The model is trained to recover x_0 by predicting each mask token x_t^i with the output probability $p_{\theta}^i(x_0^i | x_t)$. The objective is:

$$L_{\text{diffusion}}(\theta) = \mathbb{E}_{t, x_0, x_t} \left[-\frac{1}{t} \sum_{i=1}^L \mathbf{1}(x_t^i = [\text{M}]) \log p_{\theta}^i(x_0^i | x_t) \right], \quad (2)$$

where x_t is a sequence of length L , sampled from the forward process, $\mathbf{1}(\cdot)$ is an indicator function ensuring that the loss is computed only on mask tokens, following Nie et al. (2025).

During inference, the reverse diffusion process begins with a sequence of only mask tokens and proceeds for N steps at monotonically decreasing timesteps $t_{0:N}$. At step t_{n-1} , the model takes the

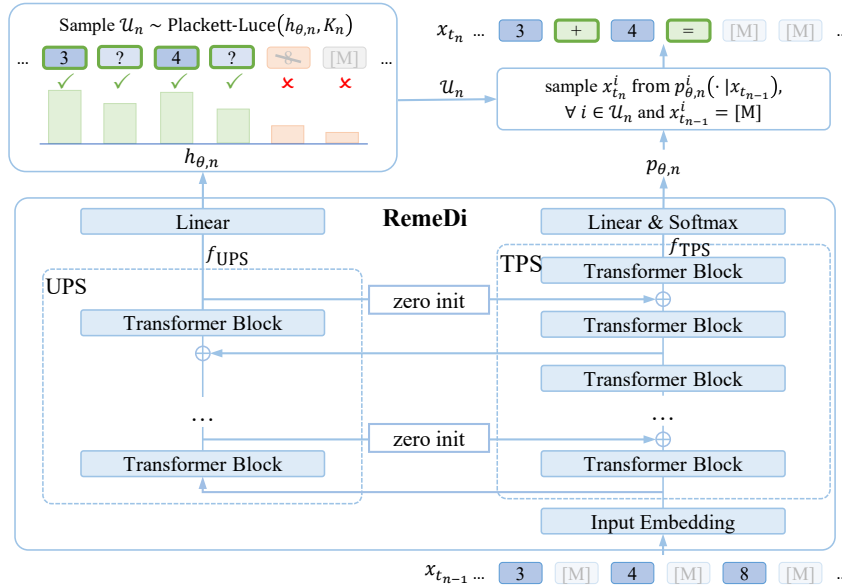


Figure 2: The structure of RemeDi, including Unmasking Policy Stream (UPS) to predict confidences h_θ for selecting the set of unmasking tokens \mathcal{U}_n , and Token Prediction Stream (TPS) to predict the token value when unmasking a masked position.

partially masked sequence $x_{t_{n-1}}$ as input and predicts all mask tokens simultaneously. A subset of tokens is unmasked to obtain x_{t_n} according to the noise schedule and the unmasking policy (e.g., unmasking tokens with the highest confidence), while the remaining predictions are remasked and deferred to later steps.

A limitation of this paradigm is that once a token is unmasked, it remains fixed in subsequent steps. In early stages, limited unmasked tokens often lead to unreliable predictions, resulting in errors that persist through the remainder of the generation process. As generation progresses, additional context may reveal these errors, but current paradigm offers no way to correct them. This motivates the ability to *remask* tokens, allowing the model to remask earlier predictions back to the mask token and predict them again using richer context in later steps.

3.2 REMEDI

We propose **RemeDi**, a DLM that can identify and remask low-confidence tokens during generation to enable iterative self-reflection. We extend the standard transformer into a dual-stream transformer architecture as shown in Fig. 2, which comprises:

- **Token Prediction Stream (TPS):** A stack of transformer blocks that predict probabilities $p_\theta^i(\cdot|x_t)$ for masked tokens as in a typical DLM (Nie et al., 2025).
- **Unmasking Policy Stream (UPS):** Another stack of transformer blocks that output token-wise confidence score h_θ^i . It represents the model’s confidence over the output tokens, indicating if they should be unmasked with high confidence. Otherwise, if the confidence is too low for a token, it should be kept masked or remasked so that it could be sampled or resampled later.

The two streams run in parallel. During inference, UPS is inserted periodically and receives hidden states from TPS as input, producing an auxiliary representation f_{UPS} for confidence scoring. The two streams perform bidirectional feature sharing: UPS layers are conditioned on f_{TPS} , and their outputs also feed back into TPS to enrich its representations. At the final layer, p and h are produced simultaneously using two independent linear heads applied to f_{TPS} and f_{UPS} , respectively. More details about the model structure for these two streams can be found in Appendix B.1.

The token generation proceeds through iterative denoising steps. Given $x_{t_{n-1}}$ as the input, UPS first predicts a confidence score $h_{\theta,n}^i$ at each position i , and select a subset of positions \mathcal{U}_n to unmask at the current step. Then, for the positions selected to be unmasked, if they have already been unmasked in $x_{t_{n-1}}$, they remain unchanged; otherwise they are sampled from $p_{\theta}^i(\cdot|x_{t_{n-1}})$ predicted by TPS. Unlike existing mask-based DLMs where tokens are fixed once being unmasked, RemeDi re-decides a token to be unmasked or (re-)masked at each step by its trained confidence score. Thus, it is possible that an already generated token is assigned with a low confidence and remasked, allowing it to be resampled in later steps. A noise schedule controls that the total number of unmasked tokens increases linearly from 0 to L (Nie et al., 2025), so that the number of mask tokens approaches zero at the final step.

In the following sections, we elaborate on how to train RemeDi with Remask SFT and Remask RL algorithms.

3.2.1 REMASK SFT

Traditional mask-based DLMs conduct SFT with randomly masked input sequences (Nie et al., 2025; Lou et al., 2024), while RemeDi needs to detect and remask possible incorrect tokens that arise during the reverse diffusion process, so they can be resampled in later steps. To achieve this, in SFT we treat such incorrect tokens as a second noise type in addition to the first noise type of mask tokens in mask-based DLMs, and train the model to recover mask tokens as well as identify unmasked tokens that should be remasked.

To simulate inference inputs at a diffusion time t , we construct training samples x_t from clean text x_0 by applying two types of noise: given a randomly sampled diffusion time $t \in (0, 1)$, we set the corresponding mask ratio $\rho_{t,\text{mask}}$, alongside the incorrect token ratio $\rho_{t,\text{incorrect}}$. With both ratios, we randomly mask tokens with $\rho_{t,\text{mask}}$. Then, among the remaining unmasked positions, we sample a subset with the ratio $\rho_{t,\text{incorrect}}$ and replace each selected token with a random alternative to simulate the incorrect tokens that may occur in the reverse diffusion process.

As aforementioned, during the reverse diffusion process, the noise level, defined as the number of mask tokens, should decrease monotonically (Guo et al., 2025). Since all incorrect tokens in an input sequence of length L must be remasked as designed below for the SFT, we require:

$$\lceil \rho_{t,\text{incorrect}} \cdot (1 - \rho_{t,\text{mask}}) \cdot L \rceil < \lceil \rho_{t,\text{mask}} \cdot L \rceil \quad (3)$$

to ensure a monotonically decreasing number of mask tokens as outputs. Otherwise, remasking all incorrect tokens would increase the total number of masks in the next step, violating the principle that the number of mask tokens should decrease at each diffusion step.

Considering the above inequality, we choose $\rho_{t,\text{mask}} = t$ and $\rho_{t,\text{incorrect}} = 4r \cdot t(1-t)$ (r is a constant) following (Nie et al., 2025; von Rütte et al., 2025). We set $r = 0.1$ in our experiments, under which it is not hard to see that the inequality 3 always holds on $t \in [0, 1]$.

Remask SFT Algorithm. During training, in addition to the typical diffusion loss in Eq. 2, we supervise the unmasking score h_{θ} with a binary cross-entropy (BCE) objective across all token positions. We construct the training label y based on different token types:

- A clean token ($i \in \mathcal{S}_{\text{clean}} = \{i \mid x_t^i = x_0^i\}$) receives a positive unmask label $y^i = 1$, indicating they should remain unmasked.
- An incorrect token ($i \in \mathcal{S}_{\text{incorrect}} = \{i \mid x_t^i \neq x_0^i, x_t^i \neq [\text{M}]\}$) receives a negative unmask label $y^i = 0$, indicating that they should be remasked.
- A mask token ($i \in \mathcal{S}_{\text{mask}} = \{i \mid x_t^i = [\text{M}]\}$) is assigned a soft unmask label $y^i = p_{\theta}^i(x_0^i|x_t)$, equal to the predicted probability of the ground-truth token x_0^i . A higher probability indicates a higher likelihood that the predicted token is correct and thus should be unmasked.

With unmask labels assigned above, we seek to minimize

$$\mathcal{L}_{\text{UPS}}(\theta) = \sum_i \text{BCE}(\sigma(h_{\theta}^i), y^i), \quad (4)$$

where $\sigma(\cdot)$ is the sigmoid function. Thus, the overall Remask SFT objective is:

$$\mathcal{L}(\theta) = \mathcal{L}_{\text{diffusion}}(\theta) + \lambda_{\text{UPS}} \mathcal{L}_{\text{UPS}}(\theta), \quad (5)$$

Algorithm 1 Input sequence construction and loss calculation in Remask SFT

Require: Clean sequence $x_0 = [x_0^1, \dots, x_0^L]$ of length L . Model \mathcal{M} with learnable parameters θ .

1: Sample $t \in (0, 1)$ according to the noise schedule, obtaining $\rho_{t,\text{mask}}$ and $\rho_{t,\text{incorrect}}$.

2: **Construct noisy input x_t :**

3: For each position i , replace x_0^i with [M] w.p. $\rho_{t,\text{mask}}$

4: Among remaining positions, replace x_0^i with a random alternative token w.p. $\rho_{t,\text{incorrect}}$

5: Define index sets:

$$\mathcal{S}_{\text{mask}} = \{i \mid x_t^i = [\text{M}]\}, \quad \mathcal{S}_{\text{incorrect}} = \{i \mid x_t^i \neq x_0^i \wedge x_t^i \neq [\text{M}]\}, \quad \mathcal{S}_{\text{clean}} = \{i \mid x_t^i = x_0^i\}$$

6: Get model outputs: $[p_\theta, h_\theta] = \mathcal{M}(x_t; \theta)$

7: Calculate the diffusion loss, on mask tokens only: $\mathcal{L}_{\text{diffusion}}(\theta) = -\frac{L}{|\mathcal{S}_{\text{mask}}|} \sum_{i \in \mathcal{S}_{\text{mask}}} \log p_\theta^i(x_0^i | x_t)$

8: Get labels for UPS: $y^i = \begin{cases} 1 & i \in \mathcal{S}_{\text{clean}} \\ 0 & i \in \mathcal{S}_{\text{incorrect}} \\ \text{stopgrad}(p_\theta^i(x_0^i | x_t)) & i \in \mathcal{S}_{\text{mask}} \end{cases}$

9: UPS BCE loss:

$\triangleright \sigma(\cdot)$ represents the sigmoid function

$$\mathcal{L}_{\text{UPS}}(\theta) = -\frac{1}{L} \sum_{i=1}^L \left(y^i \log \sigma(h_\theta^i) + (1 - y^i) \log (1 - \sigma(h_\theta^i)) \right)$$

10: Total loss: $\mathcal{L}(\theta) = \mathcal{L}_{\text{diffusion}}(\theta) + \lambda_{\text{UPS}} \mathcal{L}_{\text{UPS}}(\theta)$

where λ_{UPS} balances the two losses.

Finally, we summarize the Remask SFT in Algorithm 1, where we elaborate on how to construct the input sequence and calculate the loss function for the Remask SFT.

3.2.2 REMASK RL

After training with Remask SFT, we further fine-tune the model with outcome-based reinforcement learning (RL) to optimize the full generation trajectory (Huang et al., 2025). Specifically, we reinforce the generation process with N denoising steps, beginning from an all-mask prior $x_{t_0} \sim p_{\text{prior}}$ at $t_0 = 1$ and proceeding through timesteps $t_{0:N}$.

At each step t_n , RemeDi generates x_{t_n} from $x_{t_{n-1}}$ by invoking two coupled policies: an *unmasking policy* that chooses a subset of positions $\mathcal{U}_n = [u_n(1), \dots, u_n(K_n)]$ to unmask, and a *token prediction policy* that samples tokens at the chosen positions. Unlike standard DLMs, which never remask revealed tokens, RemeDi allows previously unmasked tokens to be remasked, enabling revision of earlier predictions.

Unmasking policy. The UPS produces a per-token confidence score h_θ^i , indicating how strongly the model believes token at position i is correct (if unmasked) or predictable (if masked). At inference, we rank tokens by their confidence scores and prioritize high-confidence ones to unmask. The number of unmasked tokens K_n at each diffusion step is determined by linearly increasing from 0 to L . During RL training, we construct an unmasking policy to sample $\mathcal{U}_n = [u_n(1), \dots, u_n(K_n)]$ using the Plackett–Luce model (Plackett, 1975): based on h_θ , we use a multinomial distribution and sequentially sample K_n positions from $\{1, \dots, L\}$ without replacement. Formally, the probability of sampling \mathcal{U}_n is:

$$\pi_{\theta,n}^{\text{unmask}}(\mathcal{U}_n \mid x_{t_{n-1}}) = \prod_{k=1}^{K_n} \frac{\exp(h_{\theta,n}^{u_n(k)})}{\sum_{j \notin \{u_n(1), \dots, u_n(k-1)\}} \exp(h_{\theta,n}^j)}. \quad (6)$$

Token prediction policy. For each position $i \in \mathcal{U}_n$, if $x_{t_{n-1}}^i = [\text{M}]$, the model samples token from $p_\theta^i(\cdot \mid x_{t_{n-1}})$; otherwise, the token remains unchanged as in the input. The probability of generating x_{t_n} given $x_{t_{n-1}}$ and \mathcal{U}_n is:

$$\pi_{\theta,n}^{\text{token}}(x_{t_n} \mid x_{t_{n-1}}, \mathcal{U}_n) = \prod_{i \in \mathcal{U}_n, x_{t_{n-1}}^i = [\text{M}]} p_\theta^i(x_{t_n}^i \mid x_{t_{n-1}}). \quad (7)$$

Joint policy. Thus, the probability of transitioning from $x_{t_{n-1}}$ to x_{t_n} is the product of the unmasking probability and the token prediction probability:

$$\pi_{\theta,n}(x_{t_n} | x_{t_{n-1}}) = \pi_{\theta,n}^{\text{unmask}}(\mathcal{U}_n | x_{t_{n-1}}) \cdot \pi_{\theta,n}^{\text{token}}(x_{t_n} | x_{t_{n-1}}, \mathcal{U}_n). \quad (8)$$

With the probability defined in Eq. 8, we apply outcome-based reinforcement learning to encourage generation trajectories $x_{t_{0:N}}$ that lead to correct final responses x_{t_N} . Specifically, we adopt GRPO (Shao et al., 2024), a scalable RL paradigm for language models. The reward is defined according to task type: verifiable correctness for math and code, and reward-model evaluation for open-ended questions. Further details on datasets and reward design are provided in Appendix B.2.4.

As shown in Fig. 11 of Appendix A, after Remask SFT and RL training, the learned h_θ serves as a reliable indicator to assess the quality of input tokens. Tokens already unmasked in the input typically receive high confidence scores. However, when certain tokens are assigned low confidence, they are more likely to be inadequate and are remasked for re-prediction in subsequent steps. It suggests that the UPS-predicted confidence scores provide a reliable estimate of per-token quality for the unmasking policy.

4 EXPERIMENTS

RemeDi enables remasking on a DLM capable of variable-length block-wise generation (Arriola et al., 2025) to support *variable-length generation*, a key feature for enabling the real-world DLM to generate an unfixed number of blocks (see Appendix B.2.2 for details). Since there are no open-source large-scale variable-length block-wise DLMs, we adapt our model from LLaDA, a widely used benchmark DLM. Starting from LLaDA’s model weights as initialization, RemeDi undergoes two stages of supervised fine-tuning and RL. We detail the training configurations in Appendix B.2, and the evaluation metrics in Appendix B.3.2.

Table 1: Model performance on math and code generation benchmarks. We highlight the best-performing model among compared DLMs in **bold**. “-” indicates unknown cases not mentioned in original papers.

Method	Math			Code	
	GSM8K	MATH	GPQA	HumanEval	MBPP
Diffusion Language Models					
Dream (Ye et al., 2025)	82.1	49.6	30.6	59.8	59.6
LLaDA (Nie et al., 2025)	78.3	38.9	28.1	45.7	39.0
LLaDA + ReMDM (Wang et al., 2025a)	81.4	38.5	-	44.5	37.8
d1-LLaDA (Zhao et al., 2025)	82.1	-	-	37.8	44.7
wd1-LLaDA (Tang et al., 2025)	82.3	-	-	-	-
LLaDA 1.5 (Zhu et al., 2025)	83.3	42.6	36.9	52.4	42.8
LLaDOU (Huang et al., 2025)	88.1	44.6	-	59.1	51.6
RemeDi (+ Remask SFT)	86.3	51.4	32.6	71.3	57.8
RemeDi (++) Remask RL)	89.1	52.9	29.5	73.2	59.4
Auto-regressive Models					
LLaMA2 7B (Touvron et al., 2023)	14.6	2.5	28.4	12.8	20.8
MetaMath 7B (Yu et al., 2023)	66.5	19.8	-	-	-
CodeLLaMA 7B (Roziere et al., 2023)	-	-	-	34.8	44.4
Deepseek 7B (Bi et al., 2024)	63.0	15.8	-	48.2	35.2
DeepseekMath 7B (Shao et al., 2024)	88.2	51.7	-	-	-
DeepseekCoder 7B (Guo et al., 2024)	-	-	-	66.1	65.4
LLaMA3 8B (Dubey et al., 2024)	78.3	29.6	31.9	59.8	57.6
Gemma2 9B (Team, 2024)	76.7	44.3	32.8	68.9	74.9

4.1 RESULTS

To evaluate various capabilities of RemeDi in different aspects, we conducted detailed comparisons against existing large language models of comparable scale in Tab. 1 and Tab. 2, including both

Table 2: Model performance on general tasks. We highlight the best-performing model among compared DLMs in **bold**. “-” indicates unknown cases not mentioned in original papers.

Method	Hellaswag	ARC-C	IFEval	AlpacaEval
Diffusion Language Models				
Dream (Ye et al., 2025)	70.3	79.2	67.5	5.9
LLaDA (Nie et al., 2025)	69.7	83.9	70.0	11.2
LLaDA 1.5 (Zhu et al., 2025)	70.5	83.5	73.5	13.9
RemeDi (+ Remask SFT)	71.1	85.2	81.9	12.5
RemeDi (++) Remask RL)	72.2	87.7	85.4	24.8
Auto-regressive Models				
LLaMA2 7B (Touvron et al., 2023)	51.5	57.3	-	-
Deepseek 7B (Bi et al., 2024)	68.5	49.4	-	-
LLaMA3 8B (Dubey et al., 2024)	75.5	82.4	-	-

DLMs and auto-regressive models. We select nine popular benchmarks across general tasks, mathematics, coding, and human preference domains.

After Remask SFT, RemeDi demonstrates high performance on almost all these benchmarks. It not only achieves the state of the art performance among existing DLMs, but also outperforms auto-regressive models of similar model size. On math benchmarks, RemeDi after Remask SFT achieves 86.3% on GSM8K and 51.4% on MATH, surpassing MetaMath with math-specific instruction tuning. It is even on par with DeepseekMath using math-specific reinforcement learning. On code generation benchmarks, RemeDi achieves 71.3% on HumanEval, outperforming CodeLLaMA and Deepseek Coder. For general natural language tasks, RemeDi also demonstrates strong performance in common knowledge answering (85.2% on ARC-C) and instruction following (81.9% on IFEval) tasks. It also aligns well with human preference (12.5% on AlpacaEval), outperforming other DLMs such as Dream and LLaDA.

After Remask RL, RemeDi achieves further improvements across a wide range of math, coding and general tasks. For example, the accuracies on GSM8K and MATH reach 89.1% and 52.9% respectively, outperforming all compared DLMs and AR models. Among all the benchmarks, RemeDi achieves its most substantial improvement on the AlpacaEval (Dubois et al., 2024) benchmark, with a +12.3% gain over the Remask SFT model. This demonstrates the effectiveness of our approach in post-training the model’s ability for a broad range of tasks.

4.2 VISUALIZATION AND ANALYSIS

We visualize how remasking improves text generation in RemeDi in Fig. 3. The model initially generated the token “making.” After generating the object “tests and estimators,” it found that “making” is not the proper verb in this verb-object structure. Thus, the model remasks it and opts for the more appropriate “developing.” This example shows RemeDi’s ability to iteratively refine its output content. We provide more examples in Appendix A, demonstrating that RemeDi is able to perform a variety of operations such as replacing, inserting and deleting with the remask mechanism.

To provide a quantitative analysis, we calculate the frequencies of remasking in a block of length 32 on MATH-500 (Lightman et al., 2023), HumanEval (Chen et al., 2021), and AlpacaEval (Dubois et al., 2024). In Fig. 4, we can see that remasking occurs most frequently in code generation, followed by mathematical reasoning, and general tasks. This pattern may be attributed to differences in structural constraints: code requires strict syntactic correctness, and mathematical solutions demand formally structured derivations, whereas responses to open-ended problems allow more flexibility.

We also analyzed remasking frequencies across different difficulty levels on MATH-500, as shown in Tab. 3. RemeDi tends to remask more frequently as the difficulty increases, rising from about 9 tokens per block at level 1–2 to nearly 14 tokens at level 4–5. This pattern suggests that iterative refinement becomes increasingly necessary for harder problems.

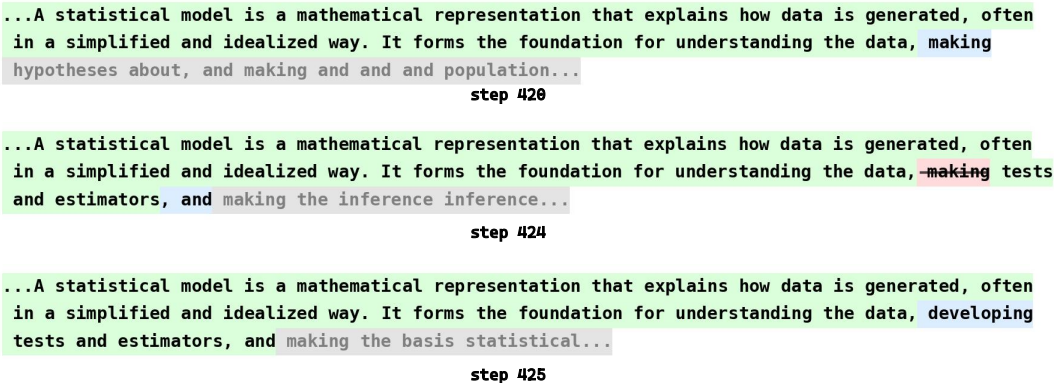


Figure 3: An example of the step-by-step generation process. Green and Red are already unmasked in the inputs. Red tokens are remarked. Blue tokens are unmasked in the outputs. Gray tokens remain masked, and we display the token with the highest probability at these positions. More examples can be found in AppendixA.

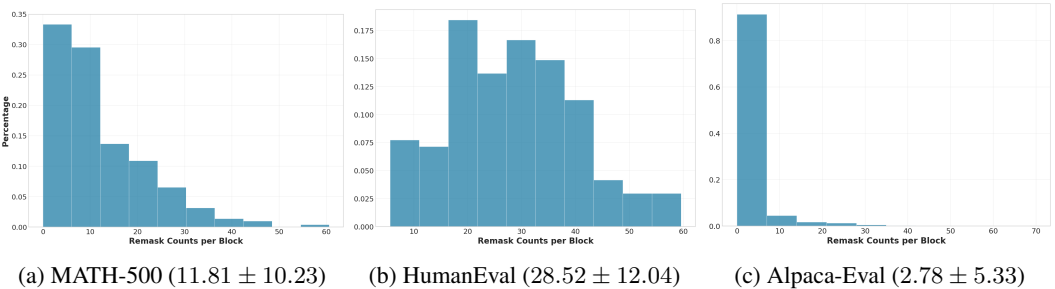


Figure 4: Distribution of remarking frequencies per block across different evaluation datasets. The numbers in parentheses indicate the mean and standard deviation for each dataset.

4.3 ABLATION STUDIES

Remark SFT We compare the improvement brought by the Remark SFT (introduced in Sec. 3.2.1) with that of vanilla SFT, under the same training configuration detailed in Appendix B.2.5. We start from a baseline model that has already completed the warm-up phase tuning for variable-length block-wise generation, and perform training on the full code-category dataset and the open2math-1M-gpt-4.1-mini dataset mentioned in Appendix B.2.1. As shown in Tab. 4, Remark SFT outperforms vanilla SFT on all benchmarks, especially on MATH-500 (+2.6%) and HumanEval (+1.8%), demonstrating that Remark SFT is an effective training method to improve DLM’s performance.

Remark RL We compare Remark RL with LLaDOU RL (Huang et al., 2025), another algorithm that also reinforces the whole generation trajectories in the reverse diffusion process. Since LLaDOU RL is developed on LLaDA, we also implement Remark RL on LLaDA for the sake of fair comparison. All experiments are conducted on GSM8K with a generation length of 256, 64 denoising steps, a block length of 64, and temperature 0.7, while all other hyperparameters follow the LLaDOU setup (see Appendix B.2.5).

Remark RL demonstrates advantages in both convergence speed and performance. As shown in Tab. 5, Remark RL achieves a higher final accuracy of 83.33%, with a particularly noticeable improvement in early training stages (e.g., 80.00% vs. 77.58% at step 50). This indicates that the more flexible remark process contributes to both faster convergence and stronger model performance.

Table 3: Statistics of the remasking frequencies per block (block size is fixed to 32) when generating responses to questions with different difficulty levels in MATH-500.

Difficulty Level	Remasking Frequencies / Block	Accuracy
1	9.13 ± 9.54	86.04%
2	8.91 ± 7.29	80.21%
3	10.13 ± 8.64	64.48%
4	13.91 ± 11.44	50.00%
5	13.95 ± 11.12	19.25%

Table 4: Experiment results after supervised tuning with different algorithms. The baseline model is already tuned to be a variable-length block-wise generation DLM (see Appendix B.2.2).

Method	GSM8K	MATH-500	HumanEval	MBPP
Baseline	80.3	34.7	41.5	42.6
Vanilla SFT	83.1	40.1	48.2	43.4
Remask SFT	83.6	42.7	50.0	44.0

Table 5: GSM8K pass@1 accuracy comparison between Remask and LLaDOU RL

Training Steps	Remask RL	LLaDOU RL
50	80.00%	77.58%
100	81.40%	78.86%
150	81.59%	80.00%
200	83.33%	82.35%

5 CONCLUSION

In this paper, we introduce the Remasking-enabled Diffusion Language Model (RemeDi), a new self-reflective remasking mechanism to address the limitation of existing mask-based DLMs that they cannot revise generated tokens. In RemeDi, remasking is achieved by predicting a confidence score to identify noisy tokens, allowing them to be remasked and then resampled with richer context in later steps.

Through a two-stage training pipeline of Remask SFT and Remask RL, RemeDi achieves the state-of-the-art performance among open-source DLMs. Our analysis further shows that the learned confidence scores provide a reliable signal of per-token quality during generation. RemeDi opens a promising direction for self-reflective text generation, further releasing the full potentials of DLMs to solve complex tasks with higher quality.

6 REPRODUCIBILITY STATEMENT

We provide a link containing the inference code and model weights, details of the datasets and configurations used in both Remask SFT and RL in Appendix B.2, and the evaluation settings in Appendix B.3.

ACKNOWLEDGEMENTS

This work was supported by the National Natural Science Foundation of China under Grant No. 92467104, and Zhejiang Leading Innovative and Entrepreneur Team Introduction Program (2024R01007).

REFERENCES

- ajibawa 2023. Maths-college. Hugging Face Datasets. URL <https://huggingface.co/datasets/ajibawa-2023/Maths-College>.
- 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.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
- Akhiad Bercovich, Itay Levy, Izik Golan, Mohammad Dabbah, Ran El-Yaniv, Omri Puny, Ido Galil, Zach Moshe, Tomer Ronen, Najeeb Nabwani, Ido Shahaf, Oren Tropp, Ehud Karpas, Ran Zilberstein, Jiaqi Zeng, Soumye Singhal, Alexander Bukharin, Yian Zhang, Tugrul Konuk, Gerald Shen, Ameya Sunil Mahabaleshwarakar, Bilal Kartal, Yoshi Suhara, Olivier Delalleau, Zijia Chen, Zhilin Wang, David Mosallanezhad, Adi Renduchintala, Haifeng Qian, Dima Rekeshe, Fei Jia, Somshubra Majumdar, Vahid Noroozi, Wasi Uddin Ahmad, Sean Narenthiran, Aleksander Ficek, Mehrzad Samadi, Jocelyn Huang, Siddhartha Jain, Igor Gitman, Ivan Moshkov, Wei Du, Shubham Toshniwal, George Armstrong, Branislav Kisacanin, Matvei Novikov, Daria Gitman, Evelina Bakhturina, Jane Polak Scowcroft, John Kamalu, Dan Su, Kezhi Kong, Markus Kliegl, Rabeeh Karimi, Ying Lin, Sanjeev Satheesh, Jupinder Parmar, Pritam Gundecha, Brandon Norrick, Joseph Jennings, Shrimai Prabhumoye, Syeda Nahida Akter, Mostofa Patwary, Abhinav Khattar, Deepak Narayanan, Roger Waleffe, Jimmy Zhang, Bor-Yiing Su, Guyue Huang, Terry Kong, Parth Chadha, Sahil Jain, Christine Harvey, Elad Segal, Jining Huang, Sergey Kashirsky, Robert McQueen, Izzy Putterman, George Lam, Arun Venkatesan, Sherry Wu, Vinh Nguyen, Manoj Kilaru, Andrew Wang, Anna Warno, Abhilash Somasamudramath, Sandip Bhaskar, Maka Dong, Nave Assaf, Shahar Mor, Omer Ullman Argov, Scot Junkin, Oleksandr Romanenko, Pedro Larroy, Monika Katariya, Marco Rovinelli, Viji Balas, Nicholas Edelman, Anahita Bhiwandiwala, Muthu Subramaniam, Smita Ithape, Karthik Ramamoorthy, Yuting Wu, Suguna Varshini Velury, Omri Almog, Joyjit Daw, Denys Fridman, Erick Galinkin, Michael Evans, Katherine Luna, Leon Derczynski, Nikki Pope, Eileen Long, Seth Schneider, Guillermo Siman, Tomasz Grzegorzec, Pablo Ribalta, Monika Katariya, Joey Conway, Trisha Saar, Ann Guan, Krzysztof Pawelec, Shyamala Prayaga, Oleksii Kuchaiev, Boris Ginsburg, Oluwatobi Olabiyi, Kari Briski, Jonathan Cohen, Bryan Catanzaro, Jonah Alben, Yonatan Geifman, Eric Chung, and Chris Alexiuk. Llama-nemotron: Efficient reasoning models, 2025. URL <https://arxiv.org/abs/2505.00949>.
- Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qishi Du, Zhe Fu, et al. Deepseek llm: Scaling open-source language models with longtermism. *arXiv preprint arXiv:2401.02954*, 2024.
- Andrew Campbell, Joe Benton, Valentin De Bortoli, Thomas Rainforth, George Deligiannidis, and Arnaud Doucet. A continuous time framework for discrete denoising models. *Advances in Neural Information Processing Systems*, 35:28266–28279, 2022.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. 2021.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *ArXiv*, abs/1803.05457, 2018.

- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems, 2021. URL <https://arxiv.org/abs/2110.14168>, 9, 2021.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pp. arXiv-2407, 2024.
- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled alpaca-eval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*, 2024.
- Nima Fathi, Torsten Scholak, and Pierre-André Noël. Unifying autoregressive and diffusion-based sequence generation. *arXiv preprint arXiv:2504.06416*, 2025.
- Itai Gat, Heli Ben-Hamu, Marton Havasi, Daniel Haziza, Jeremy Reizenstein, Gabriel Synnaeve, David Lopez-Paz, Brian Karrer, and Yaron Lipman. Set block decoding is a language model inference accelerator. *arXiv preprint arXiv:2509.04185*, 2025.
- Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming—the rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024.
- Zhehao Guo, Jiedong Lang, Shuyu Huang, Yunfei Gao, and Xintong Ding. A comprehensive review on noise control of diffusion model. *arXiv preprint arXiv:2502.04669*, 2025.
- Marton Havasi, Brian Karrer, Itai Gat, and Ricky TQ Chen. Edit flows: Flow matching with edit operations. *arXiv preprint arXiv:2506.09018*, 2025.
- Zhiwei He, Tian Liang, Jiahao Xu, Qiuzhi Liu, Xingyu Chen, Yue Wang, Linfeng Song, Dian Yu, Zhenwen Liang, Wenxuan Wang, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. Deepmath-103k: A large-scale, challenging, decontaminated, and verifiable mathematical dataset for advancing reasoning. 2025. URL <https://arxiv.org/abs/2504.11456>.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021.
- Chihan Huang and Hao Tang. CtrlDiff: Boosting large diffusion language models with dynamic block prediction and controllable generation. *arXiv preprint arXiv:2505.14455*, 2025.
- Siming Huang, Tianhao Cheng, Jason Klein Liu, Jiaran Hao, Liuyihan Song, Yang Xu, J. Yang, J. H. Liu, Chenchen Zhang, Linzheng Chai, Ruifeng Yuan, Zhaoxiang Zhang, Jie Fu, Qian Liu, Ge Zhang, Zili Wang, Yuan Qi, Yinghui Xu, and Wei Chu. Opencoder: The open cookbook for top-tier code large language models. 2024. URL <https://arxiv.org/pdf/2411.04905>.
- Zemin Huang, Zhiyang Chen, Zijun Wang, Tiancheng Li, and Guo-Jun Qi. Reinforcing the diffusion chain of lateral thought with diffusion language models. *arXiv preprint arXiv:2505.10446*, 2025.
- Jaeyeon Kim, Kulin Shah, Vasilis Kontonis, Sham Kakade, and Sitan Chen. Train for the worst, plan for the best: Understanding token ordering in masked diffusions. *arXiv preprint arXiv:2502.06768*, 2025.
- Hynek Kydlíček. Math-verify: A robust mathematical expression evaluation system. <https://github.com/huggingface/Math-Verify>, 2025.
- Jia LI, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Costa Huang, Kashif Rasul, Longhui Yu, Albert Jiang, Ziju Shen, Zihan Qin, Bin Dong, Li Zhou, Yann Fleureau, Guillaume Lample, and Stanislas Polu. NuminaMath. [<https://huggingface.co/AI-MO/NuminaMath-CoT>] (https://github.com/project-umina/aimo-progress-prize/blob/main/report/umina_dataset.pdf), 2024.
- Jijie Li, Li Du, Hanyu Zhao, Bowen Zhang, Liangdong Wang, Boyan Gao, Guang Liu, and Yonghua Lin. Infinity instruct: Scaling instruction selection and synthesis to enhance language models, 2025. URL <https://arxiv.org/abs/2506.11116>.

- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*, 2023.
- Chris Yuhao Liu, Liang Zeng, Jiakai Liu, Rui Yan, Jujie He, Chaojie Wang, Shuicheng Yan, Yang Liu, and Yahui Zhou. Skywork-reward: Bag of tricks for reward modeling in llms. *arXiv preprint arXiv:2410.18451*, 2024.
- Chris Yuhao Liu, Liang Zeng, Yuzhen Xiao, Jujie He, Jiakai Liu, Chaojie Wang, Rui Yan, Wei Shen, Fuxiang Zhang, Jiacheng Xu, Yang Liu, and Yahui Zhou. Skywork-reward-v2: Scaling preference data curation via human-ai synergy. *arXiv preprint arXiv:2507.01352*, 2025.
- Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion modeling by estimating the ratios of the data distribution. In *Forty-first International Conference on Machine Learning*, 2024.
- Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Tianjun Zhang, Erran Li, Raluca Ada Popa, and Ion Stoica. Deepscaler: Surpassing o1-preview with a 1.5b model by scaling rl, 2025. Notion Blog.
- Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with evol-instruct, 2023.
- Konstanty Marczak. science_qa. Hugging Face Datasets, 2023. URL https://huggingface.co/datasets/KonstantyM/science_qa.
- Arindam Mitra, Hamed Khanpour, Corby Rosset, and Ahmed Awadallah. Orca-math: Unlocking the potential of slms in grade school math, 2024.
- mlfoundations dev. open2math-1m-gpt-4.1-mini. Hugging Face Datasets, 2025. URL <https://huggingface.co/datasets/mlfoundations-dev/open2math-1M-gpt-4.1-mini>.
- Dhruv Nathawani, Shuoyang Ding, Vitaly Lavrukhin, Igor Gitman, Somshubra Majumdar, Evelina Bakhturina, Boris Ginsburg, and Jane Polak Scowcroft. Nemotron-Post-Training-Dataset-v2, August 2025. URL <https://huggingface.co/datasets/nvidia/Nemotron-Post-Training-Dataset-v2>.
- 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.
- 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.
- Robin L Plackett. The analysis of permutations. *Journal of the Royal Statistical Society Series C: Applied Statistics*, 24(2):193–202, 1975.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark, 2023.
- Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.
- 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.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.

- 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.
- Xiaohang Tang, Rares Dolga, Sangwoong Yoon, and Ilija Bogunovic. wd1: Weighted policy optimization for reasoning in diffusion language models. *arXiv preprint arXiv:2507.08838*, 2025.
- Gemma Team. Gemma. 2024. doi: 10.34740/KAGGLE/M/3301. URL <https://www.kaggle.com/m/3301>.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutvi Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Dimitri von Rütten, Janis Fluri, Yuhui Ding, Antonio Orvieto, Bernhard Schölkopf, and Thomas Hofmann. Generalized interpolating discrete diffusion. *arXiv preprint arXiv:2503.04482*, 2025.
- Guanghan Wang, Yair Schiff, Subham Sekhar Sahoo, and Volodymyr Kuleshov. Remasking discrete diffusion models with inference-time scaling. *arXiv preprint arXiv:2503.00307*, 2025a.
- 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*, 2025b.
- Yunhui Xia, Wei Shen, Yan Wang, Jason Klein Liu, Huifeng Sun, Siyue Wu, Jian Hu, and Xiaolong Xu. Leetcodedataset: A temporal dataset for robust evaluation and efficient training of code llms. *arXiv preprint arXiv:2504.14655*, 2025.
- Tengyu Xu, Eryk Helenowski, Karthik Abinav Sankararaman, Di Jin, Kaiyan Peng, Eric Han, Shao-liang Nie, Chen Zhu, Hejia Zhang, Wenxuan Zhou, et al. The perfect blend: Redefining rlhf with mixture of judges. *arXiv preprint arXiv:2409.20370*, 2024.
- Zhangchen Xu, Yang Liu, Yueqin Yin, Mingyuan Zhou, and Radha Poovendran. Kodcode: A diverse, challenging, and verifiable synthetic dataset for coding. *arXiv preprint arXiv:2503.02951*, 2025.
- Jiacheng Ye, Zhihui Xie, Lin Zheng, Jiahui Gao, Zirui Wu, Xin Jiang, Zhenguo Li, and Lingpeng Kong. Dream 7b, 2025. URL <https://hkunlp.github.io/blog/2025/dream>.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*, 2023.
- Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. Mammoth: Building math generalist models through hybrid instruction tuning. *arXiv preprint arXiv:2309.05653*, 2023.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019.
- Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 3836–3847, 2023.
- Siyao 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.
- Lin Zheng, Jianbo Yuan, Lei Yu, and Lingpeng Kong. A reparameterized discrete diffusion model for text generation. *arXiv preprint arXiv:2302.05737*, 2023.
- Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models, 2023. URL <https://arxiv.org/abs/2311.07911>.
- 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.

A GENERATION PROCESS OF REMEDI

To better understand how RemeDi leverages the remasking mechanism, we visualize intermediate steps when solving math, code, and open-ended problems. Since the full responses are usually long, we focus on the token segment where critical remasking occurs. In the following figure, **green** boxes indicate already generated tokens, **blue** boxes represent tokens unmasked in this step, **red** boxes denote tokens remasked in this step, and **gray** boxes represent tokens that remain masked, showing the token with the highest probability. Key tokens are highlighted with bounding boxes.

We find that remasking enables diverse forms of revision beyond simple correction, including:

- **Correcting calculation errors:** Remask can correct calculation errors. As shown in Fig. 5, the model initially predicted “\div” as the most probable operator and generated “0”. However, since the actual operator generated was “\mod”, the model remasked the previous “0” and regenerated “5” as the correct result.
- **Refining text quality:** Remasking allows more precise wording. In Fig 6, the initially generated phrase “methyl group” is not adequate when answering this problem. RemeDi replaces them with more precise term “secondary carbon” by remasking.
- **Merging adjacent tokens:** When two consecutive tokens correspond to a single vocabulary token, RemeDi may remask them and merge into one, thereby freeing a slot for subsequent generation. In Fig. 7, the separate tokens “,” and “\” were remasked and merged into the single token “,\”, releasing one token slot.
- **Splitting tokens:** Conversely, the model can split a token into smaller parts to fill idle positions, ensuring that no tokens remain unused and the denoising process can complete. In Fig. 8, to fill in the slot before “Mish”, RemeDi remasks “Mish” and regenerates it as two tokens, “M” and “ish”.
- **Inserting tokens:** Remasking also supports insertion. In Fig. 9, to add the word “again” before “bounces”, the model first remasked the two tokens “b” and “ounces”, and then regenerated the sequence with the insertion.
- **Deleting tokens:** Finally, remasking can delete tokens and replace them with nothing or control symbols. In Fig. 10, the phrase “per hour →” was removed and replaced with a line break.

These cases illustrate that remasking gives RemeDi considerable freedom to revise its outputs in multiple ways, greatly extending the flexibility of diffusion-based text generation.

Moreover, we illustrate the predicted confidence scores in Fig. 11. In general, unmasked tokens receive higher confidence scores h_θ than masked tokens, unless the model judges them as unreliable and decides to remask them. For example, see the tokens “say” in Fig. 11b and “in” in Fig. 11c. Interestingly, these tokens that are eventually remasked already exhibit relatively low confidence at the step when they were first predicted, as reflected by the lighter background green shading. This suggests that the model was uncertain about them from the start. Once more context is revealed in subsequent steps, RemeDi is able to revise such low-confidence tokens into more appropriate alternatives.

B EXPERIMENT DETAILS

B.1 DUAL-STREAM MODEL STRUCTURE

We construct TPS with the same transformer structure as LLaDA (Nie et al., 2025), comprising 32 transformer blocks. The model weights in this stream are initialized with LLaDA-8B-Instruct. For the UPS, we stack four transformer blocks with the same hidden dimension to construct a smaller network with random initialization. These two streams are weakly coupled via bi-directional connections at TPS blocks 1, 11, 21, and 31. At each connection point, the output from the previous TPS block is added to the current UPS feature to form the input for the next UPS block, while the output of the current UPS block is added to the output of the corresponding TPS block before it is passed onward. Thus, both streams enrich their representation with features from each other. To preserve the TPS’s original capability inherited from the pretrained weights, we add a zero-initialized



Figure 5: An example of correcting calculation errors with remasking. Question: A group of N students, where $N < 50$, is on a field trip. If their teacher puts them in groups of 8, the last group has 5 students. If their teacher instead puts them in groups of 6, the last group has 3 students. What is the sum of all possible values of N ?



Figure 6: An example of refining text quality with remasking. Question: What is the major outcome of the reaction between 4,4-dimethylcyclopent-1-enol and bromine?



Figure 7: An example of merging adjacent tokens. Question: *A quantum mechanical particle of mass m moves in two dimensions in the following potential, as a function of $(r, \theta) : V(r, \theta) = 1/2kr^2 + 3/2kr^2 \cos^2(\theta)$ Find the energy spectrum.*



Figure 8: An example of splitting tokens. Question: *Mishka bought 3 pairs of shorts, 3 pairs of pants, and 3 pairs of shoes. One pair of shorts costs \$16.50. One pair of pants costs \$22.50 and one pair of shoes costs \$42. How many dollars did Mishka spend on all the clothing items?*



Figure 9: An example of inserting tokens. Question: *Nathan has a bouncy ball that bounces to $2/3$ rd's of its starting height with each bounce. If he drops it from the third-floor balcony in the mall, where each story is 24 feet high, how high does the ball go on its second bounce?*

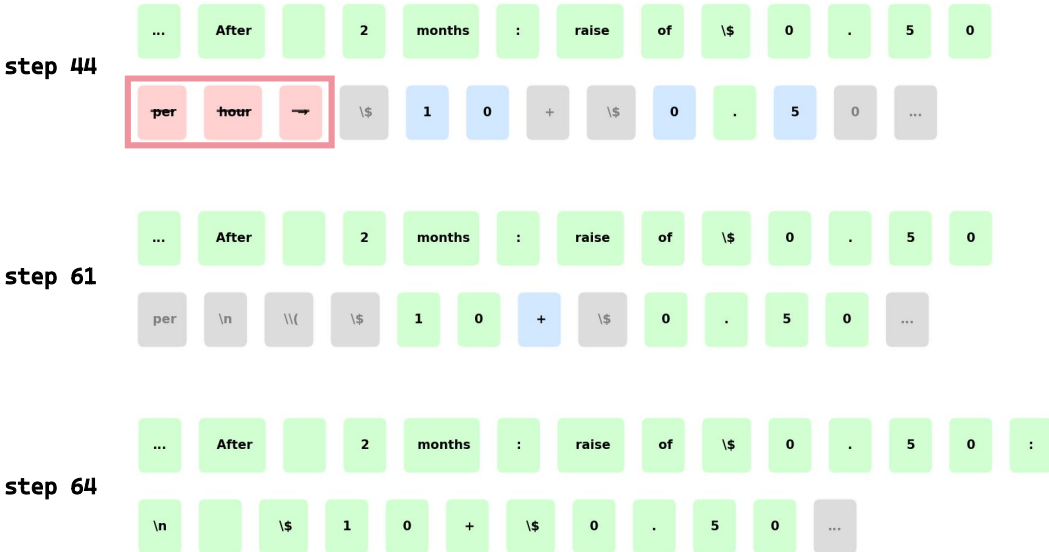


Figure 10: An example of deleting tokens. Question: *When Billy was first hired, he was paid at a rate of \$10 per hour. After 2 months, he was given a raise of \$0.50 per hour. On his first anniversary at work, he was given a raise of \$1.00 per hour. Sally just started working at a different business, and her starting salary is \$0.50 more per hour than Billy’s starting salary was. If both Billy and Sally work 20 hours, how much more money will Billy earn than Sally, in dollars?*

projection (Zhang et al., 2023) on the connections from UPS to TPS. This ensures that the model’s token prediction behavior is unchanged at the beginning of training, and gradually learns to predict the confidence for unmasking tokens in a diffusion process. The model comprises a total of 8.9B parameters.

B.2 TRAINING CONFIGURATIONS AND DATASETS

B.2.1 DATASETS

We use both high quality public datasets and in-house data for training, including four major categories: mathematics, code, general conversation, and science. The mathematics category includes NuminaMath-CoT (LI et al., 2024), MetaMathQA (Yu et al., 2023), orca-math-word-problems-200k (Mitra et al., 2024), Maths-College (ajibawa 2023), DeepMath-103K (He et al., 2025), MathInstruct (Yue et al., 2023), and open2math-1M-gpt-4.1-mini (mlfoundations dev, 2025). The code category comprises evol-codealpaca-v1 (Luo et al., 2023), opc-sft-stage1 (Huang et al., 2024), and KodCode-V1-SFT-4o (Xu et al., 2025). General conversation data is drawn from open-perfectblend (Xu et al., 2024) and Infinity-Instruct (Li et al., 2025). The science category incorporates the science_qa (Marczak, 2023) dataset. In total, the public datasets provide roughly 18.8M samples. The in-house data, containing about 140K samples of prompt-response pairs, was primarily generated by GPT-4.1 and covers mathematics, code, science, and instruction-following tasks. These generated samples went through a human check to ensure overall quality.

B.2.2 VARIABLE-LENGTH BLOCK-WISE GENERATION

We first adapt LLaDA into a DLM capable of variable-length generation. The underlying architecture remains unchanged, but we fine-tune it for variable-length block-wise generation. During inference, generation proceeds block by block, where each block consists of $L = 32$ tokens. For each block, the model runs a full reverse diffusion process until the block is fully denoised, after which the completed block is appended to the context. The next block is then generated in the same manner, continuing until an `<eos>` token appears. Similar to auto-regressive (AR) models, we enforce causality using a block-wise causal mask in the self-attention layers: each token can attend to

all tokens within its current block and all tokens in previously generated blocks, but never to tokens in future blocks.

For supervised finetuning, the response part is divided into blocks of length $L = 32$, with the last incomplete block padded by $\langle \text{eos} \rangle$ tokens. The training objective is to recover a noised version of each block conditioned on all previous clean blocks. Following Arriola et al. (2025), we concatenate both the clean blocks and their corresponding noised versions into a single input sequence, allowing all blocks to be trained jointly in one forward-backward pass.

In this stage, we train variable-length block-wise generation with learning rate 2×10^{-6} , batch size 160, and a gradient threshold of 1.0. The baseline model in Table 4 reports the results of the variable-length block-wise model on several datasets, alongside the comparisons when Remask SFT and Remask RL added to this baseline model.

B.2.3 REMASK SFT

After finetuning LLaDA for variable-length block-wise generation, we attach the Unmasking Policy Stream (UPS) to construct the model architecture shown in Fig. 2. We then further train RemeDi with Remask SFT with $\lambda_{UPS} = 0.3$. For optimization, we apply a learning rate of 2.0×10^{-5} to the newly introduced parameters in UPS, while keeping the learning rate for the original parameters at 2.0×10^{-6} .

B.2.4 REMASK RL

To enable effective RL training, we curated a dataset spanning mathematics, coding, STEM, instruction-following and preference-alignment tasks:

- **Math:** GSM8K(Cobbe et al., 2021), MATH(Hendrycks et al., 2021), DeepScaleR(Luo et al., 2025)
- **Code:** KodCode-V1-SFT-R1(Xu et al., 2025), LeetCodeDataset(Xia et al., 2025)
- **General:** Skywork-Reward-Preference-80K-v0.2(Liu et al., 2024), Llama-Nemotron-Post-Training-Dataset-RL (instruct-following)(Bercovich et al., 2025), Nemotron-Post-Training-Dataset-v2 (stem)(Nathawani et al., 2025)

To ensure quality, we applied the following filters:

- **Length:** Since our RL training limits generation length to 1024 tokens, we discard any sample—question plus answer or question alone—exceeding this bound.
- **Verifiable:** For math data, When both a short answer and a detailed response are available, we keep the sample only if the two answers match; for code data, we require that the provided solution passes all test cases.
- **Deduplication:** Given the diverse sources, we perform global deduplication using Min-HashLSH.

Reward Design. Our reward function incorporates two distinct types of reward signals:

- **Verifiable Reward:** Verifiable rewards are widely used in mathematics, code, and STEM domains, where the answer is first extracted and then verifies it: math and STEM via Math-Verify (Kydlicek, 2025), and code via executing test cases and computing the pass rate. We also incorporated verifiable instruction-following samples with IFEval (Zhou et al., 2023) format to further improve the model’s ability to follow instructions.
- **Model-based Reward:** We incorporated the Skywork-Reward-V2-Llama-3.1-8B (Liu et al., 2025), which was trained on human preference data, to evaluate response quality. Each response is assigned a scalar reward score reflecting human preference, thereby enhancing the model’s ability to produce outputs that better align with human preference during RL training.

We optimized the model using the AdamW optimizer with a learning rate of 5.0×10^{-6} , $\beta = (0.9, 0.999)$ and a maximum gradient norm of 1.0, for a total of 100 training steps.

Algorithm 2 Remask RL

Require: Model parameters θ , a dataset \mathcal{D} , and reward_func.

- 1: **while** θ not converged and maximum epochs not reached **do**
- 2: Sample questions $q \sim \mathcal{D}$
- 3: **for** $g = 1$ **to** G **do** ▷ Generate a group of G trajectories
- 4: Initialize $x_{t_0}^g$ with q and mask tokens.
- 5: **for** $n = 1$ **to** N **do** ▷ N denotes the number of denoising steps
- 6: Calculate the ranking score $h_{\theta,n}$ for each token
- 7: Sample K_n positions to unmask in this step: $\mathcal{U}_n \sim \text{Plackett-Luce}(h_{\theta,n}, K_n)$
- 8: Sample $x_{t_n}^{g,i} \sim p_{\theta,n}^i(\cdot | x_{t_{n-1}}^g)$, $\forall i \in \mathcal{U}_n, x_{t_{n-1}}^{g,i} = [M]$
- 9: **end for**
- 10: $r^g = \text{reward_func}(q, x_{t_N}^g)$ ▷ Compute the rewards
- 11: **end for**
- 12: **for** $g = 1$ **to** G **do** ▷ Compute the advantages as in GRPO
- 13: $A^g = \frac{r^g - \text{mean}(r^{1:G})}{\text{std}(r^{1:G})}$
- 14: **end for**
- 15: **for** $n = 1$ **to** N **do** ▷ Compute π_θ and losses for each denoising step
- 16: $\pi_{\theta,n}(x_{t_n}^g | x_{t_{n-1}}^g) = \pi_{\theta,n}^{\text{unmask}}(\mathcal{U}_n^g | x_{t_{n-1}}^g) \cdot \pi_{\theta,n}^{\text{token}}(x_{t_n}^g | x_{t_{n-1}}^g, \mathcal{U}_n^g)$ ▷ see Eq. 8
- 17: $\mathcal{L}_{\theta,n} = -\frac{1}{G} \sum_{g=1}^G \frac{\pi_{\theta,n}(x_{t_n}^g | x_{t_{n-1}}^g)}{\pi_{\text{old},n}(x_{t_n}^g | x_{t_{n-1}}^g)} A^g$
- 18: Calculate the gradient $\nabla_\theta \mathcal{L}_{\theta,n}$
- 19: **end for**
- 20: Update θ with accumulated gradients $\sum_{n=1}^N \nabla_\theta \mathcal{L}_{\theta,n}$ along the descent direction
- 21: **end while**

B.2.5 ABLATION SETUP

We provide detailed configurations for ablation studies in Sec. 4.3.

Remask SFT Both the vanilla SFT model and the Remask SFT model are trained with identical hyper-parameters: block size 32, maximum generation length 1024, global batch size 80, and a total of 3000 training steps. We use the AdamW optimizer with a learning rate of 2.0×10^{-5} in newly added parameters of UPS and of 2.0×10^{-6} in original TPS parameters, $\beta = (0.9, 0.999)$, and a maximum gradient norm of 1.

Remask RL Both Remask RL and LLaDOU RL are trained on LLaDA using identical RL hyper-parameters. Specifically, roll-outs are generated with temperature 0.7, generation length 256, block length 64, and $N = 64$ denoising steps. Each batch consists of 16 prompts, with each prompt generating $G = 16$ roll-outs, for a total of 200 training steps. We use the AdamW optimizer with a learning rate of 5.0×10^{-6} , $\beta = (0.9, 0.999)$, and a maximum gradient norm of 1. The complete training algorithm is elaborated in Alg. 2.

B.3 EVALUATION DETAILS**B.3.1** INFERENCE SETTINGS

Remedi For evaluation, Remedi uses a maximum generation length of 2048 on MATH and 1024 on all other datasets. At each step, only one token is unmasked, with a block size of 32 for generation. Both TPS and UPS adopt greedy sampling.

LLaDA The evaluation of LLaDA largely follows (Nie et al., 2025). On GSM8K and MATH, we set the generation length to 256 with a block length of 8, unmasking one token per step in a semi-autoregressive manner with greedy sampling. On HumanEval and MBPP, we use a generation length of 512 with a block length of 32, while keeping all other settings unchanged.

LLaDA + ReMDM We implemented ReMDM (Wang et al., 2025a) on top of LLaDA (Nie et al., 2025). Specifically, we adopted the ‘‘ReMDM-cap + Switch’’ configuration with $\eta_{\text{cap}} = 0.04$ and

$t_{\text{switch}} = 0.55$. For evaluation, we set the generation length to 256/256/512/512 and the block length to 8/8/32/32 for GSM8K, MATH, HumanEval, and MBPP, respectively.

B.3.2 BENCHMARKS

Here we provide the detailed input prompts and how the metrics are computed for different benchmarks:

GSM8K GSM8K evaluates multi-step mathematical reasoning in elementary problems (Cobbe et al., 2021). We illustrate below a zero-shot prompt used to evaluate the model. After generation, we extract the answer in “boxed{ }” from the response, and check if it is equivalent to the ground truth with the scripts developed by Hendrycks et al. (2021). We report the accuracy on this benchmark.

Janet’s ducks lay 16 eggs per day. She eats three for breakfast every morning and bakes muffins for her friends every day with four. She sells the remainder at the farmers’ market daily for \$2 per fresh duck egg. How much in dollars does she make every day at the farmers’ market? (Please put the final answer in \boxed{} tag, i.e. \$\boxed{\text{answer here}}\$)

MATH MATH contains 5,000 challenging competition mathematics problems (Hendrycks et al., 2021). We evaluate the model in a zero-shot setting, with prompts like the one below. After generation, we extract the answer in “boxed{ }” from the response, and verify if it is equivalent to the ground truth with the scripts developed by Hendrycks et al. (2021).

Convert the point $(0,3)$ in rectangular coordinates to polar coordinates. Enter your answer in the form (r,θ) , where $r > 0$ and $0 \leq \theta < 2\pi$. (Please put the final answer in \boxed{} tag, i.e. \$\boxed{\text{answer here}}\$)

GPQA GPQA is a challenging multiple-choice benchmark for testing LLM’s complex scientific reasoning and specialized knowledge domains (Rein et al., 2023). We used all 448 questions from the main version of GPQA and evaluated the model in a zero-shot setting with the prompt shown below. We select the token with the highest probability at the <mdm_mask> position as the final answer, and report the pass@1 on this benchmark.

```
<|startoftext|><|start_header_id|>user<|end_header_id|>
```

What is the correct answer to this question: A large gene has dozens of exons, of which the central ones code for folded triple helical repeats that connect the cytoskeleton with sarcolemma and extracellular space. Each exon usually codes for one folded triple alpha helix. The most common mutations of the gene are central exon deletions that create out-of-frame peptides and progressive degenerative organ waste. A solution is to deliver a Morpholino that recognizes the 5’ end of the out-of-frame exon in pre-mRNA. The molecule prevents binding of the spliceosome and creates exon skipping and in-frame joining. Several missing exons are well tolerated by an organism. Which structure below is not involved in the proposed therapy? Choices:

- (A) lariat
- (B) R-loops
- (C) antisense
- (D) polyA tail

Answer:Your answer should in the format ‘The best answer is [the_answer_letter]’ where the [the_answer_letter] is one of (A), (B), (C) or (D).<|eot_id|><|start_header_id|>assistant<|end_header_id|>

The best answer is <mdm_mask>.

MBPP MBPP consists of 500 python programming problems for entry level programmers (Austin et al., 2021). We evaluate the model with the prompt below in a zero-shot setting. After generation, we extract the python code from the response, and check if it passes all test cases associated with this problem, and report pass@1 on this benchmark.

You are an expert Python programmer. Your task is to complete the implementation of a function named `remove_Occ`.

**** TARGET FUNCTION ****

Write a python function to remove first and last occurrence of a given character from the string.

**** UNIT TESTS ****

Your code should pass unit tests like:

```
assert remove_Occ("hello", "l") == "heo"
assert remove_Occ("abcda", "a") == "bcd"
assert remove_Occ("PHP", "P") == "H"
```

Here is the function to complete:

```
```python
def remove_Occ(input_param_1, input_param_2):
 """Write a python function to remove first and last occurrence of a
 given character from the string."""
 ...
```

**HumanEval** HumanEval consists of 164 hand-written programming problems (Chen et al., 2021). We evaluate the model on it with the prompt below in a zero-shot setting. After generation, we extract the Python code from the response, and check whether the output function passes all the provided test cases; we then report the pass@1 on this benchmark.

You are an expert Python programmer, Python code should be placed between the line of ```python and the line of ``` for easy extraction later, and here is your task:

```
```python
from typing import List
def has_close_elements(numbers: List[float], threshold: float) -> bool:
    """Check if in given list of numbers, any two numbers are closer to
    each
    other than the given threshold.
```

Examples:

```
>>> has_close_elements([1.0, 2.0, 3.0], 0.5)
False
>>> has_close_elements([1.0, 2.8, 3.0, 4.0, 5.0, 2.0], 0.3)
True
"""
...

```

Hellaswag Hellaswag is a benchmark dataset specifically designed to evaluate machine common-sense reasoning capabilities (Zellers et al., 2019). It primarily assesses a model’s ability to infer the most plausible subsequent event based on given contextual information in natural language understanding tasks. We evaluate the model in a zero-shot setting on Hellaswag. We follow the approach of (Nie et al., 2025), incorporating Classifier-Free Guidance (CFG) and set the CFG weight to 0.5. Under CFG intervention, the model simultaneously computes conditional predictions (based on the given context) and unconditional predictions (absent specific context), guiding the generation process by scaling the difference between them. The final accuracy is calculated based on the model’s normalized probability assigned to the correct option.

ARC-C ARC-C is a highly challenging benchmark dataset specifically designed to evaluate machine abstract reasoning and scientific problem-solving capabilities (Clark et al., 2018). We evaluate the model in a zero-shot setting on ARC-C, with the prompt below. After generation, we extract the answer after ‘The best answer is’, and report the pass@1 rate on this benchmark.

Given the following question and four candidate answers (A, B, C and D), choose the best answer.
 Question: An astronomer observes that a planet rotates faster after a meteorite impact. Which is the most likely effect of this increase in rotation?
 A. Planetary density will decrease.
 B. Planetary years will become longer.
 C. Planetary days will become shorter.
 D. Planetary gravity will become stronger.
 Your response should end with "The best answer is [the_answer_letter]" where the [the_answer_letter] is one of A, B, C or D.

IFEval IFEval evaluates the model’s instruction-following capability with verifiable instructions (Zhou et al., 2023). We use the official evaluation code provided by the IFEval Benchmark, and compute the model’s accuracy based on the loose metric.

AlpacaEval AlpacaEval evaluates the model’s instruction-following capability with the LLM-as-a-Judge methodology (Dubois et al., 2024). As officially recommended by the AlpacaEval benchmark, we use GPT4-1106-preview as the baseline/reference model and the weighted_alpaca_eval_gpt4_turbo as the evaluator/annotator, and assess the win rate of the responses generated by RemeDi under length-controlled conditions to eliminate the confounding effect of response length.

C MORE EXPERIMENTS

C.1 COMPARISON WITH SEED DIFFUSION

Since the official Seed Diffusion model and implementation are not publicly available, we do our best to reproduce it. For a fair comparison, we train seed diffusion under the same base model and identical training configuration as in Appendix B.2.5. As shown in Table 6, Remask SFT consistently outperforms Seed Diffusion, demonstrating the advantage of learning an explicit remasking policy during training.

Table 6: Unified head-to-head comparison with other training algorithms under identical settings.

Method	GSM8K	MATH-500	HumanEval	MBPP
Baseline	80.3	34.7	41.5	42.6
Vanilla SFT	83.1	40.1	48.2	43.4
Seed Diffusion	63.9	28.0	5.4	9.8
Remask SFT	83.6	42.7	50.0	44.0

C.2 PREDICTOR-CORRECTOR VS. LEARNED REMASK POLICY

We apply a representative predictor-corrector sampler, ReMDM (Wang et al., 2025a), to RemeDi-Instruct. The evaluation setup is the same as in Appendix B.3.1. The results in Table 7 show that our learned remasking policy (via Remask SFT) is more effective than the random remasking strategy employed in predictor-corrector samplers.

Table 7: Comparison between our learned remask policy in RemeDi and the ReMDM predictor-corrector, both evaluated with RemeDi-Instruct.

Method	GSM8K	MATH-500	HumanEval	MBPP
RemeDi + predictor-corrector	58.3	38.7	39.6	54.2
RemeDi (Ours)	86.3	52.2	71.3	57.8

C.3 EFFECT OF DIFFERENT SAMPLERS

To isolate the effect of our multi-task objective and incorrect-token augmentation, we compared vanilla sampler, adaptive sampler (Kim et al., 2025) and our remask sampler under the same Remask SFT model in Table 8.

Table 8: Effect of different samplers under the same Remask SFT model.

Sampler	GSM8K	MATH-500	HumanEval	MBPP	IFEval
Vanilla	86.3	38.3	43.9	55.4	69.2
Adaptive (Kim et al., 2025)	86.6	40.3	43.9	55.4	69.2
Remask (Ours)	86.3	52.2	71.3	57.8	81.9

C.4 MATCHED-COMPUTE ABLATION: EXTRA SFT VS. REMASK RL

To directly address whether performance gains come from RL or merely additional training time, we conduct a matched-compute ablation: we continue training the RemeDi-Instruct model for an extra 2,000 steps (which consumes approximately 32 H800-days, the same compute as used in the RL stage) and evaluate its performance. The results are summarized in Table 9:

Table 9: Matched-compute ablation between extra SFT training and Remask RL.

	GSM8K	MATH-500	HumanEval	MBPP
RemeDi-Instruct	86.3	52.2	71.3	57.8
+ ~32 H800-days SFT training	83.6	52.6	62.8	57.8
+ ~32 H800-days RL training	89.1	53.2	73.2	59.4

C.5 UPS STRUCTURE ABLATIONS

Table 10 reports ablations on UPS components, showing that removing either the bi-residual connections or the zero-init bridge leads to clear performance degradation. We train all models under the same base model and identical training configuration as in Appendix B.2.5

Table 10: UPS structure ablations. Removing either the bi-residual connections or the zero-init bridge degrades performance

	GSM8K	MATH-500	HumanEval	MBPP
Baseline	83.6	42.7	50.0	44.0
w/o bi-residuals	76.6	43.1	45.7	42.2
w/o zero-init	75.2	42.7	45.1	44.8

C.6 EFFECTIVENESS OF RL

To assess the effect of RL, we compare the per-token remask frequencies between the SFT-only model and the RL-trained model. Interestingly, the RL-trained model performs **more** remasking on average. As shown in Table 11, higher remask frequency correlates with improved performance, suggesting that additional remasking provides more opportunities for RemeDi to detect and correct early-step errors.

D EXTRA VISUALIZATION

Fig. 12 compares the per-token remask frequency of the SFT-only model and the RL-trained model on GSM8K. We observe that RL training consistently increases the remask frequency, suggesting

Table 11: Average remask frequency (ARF) and performance across tasks. ARF measures how many times each token is remasked on average during decoding.

Model	GSM8K		HumanEval		AlpacaEval	
	Acc	ARF	Acc	ARF	Acc	ARF
RemeDi (+ Remask SFT)	86.3	0.16	71.3	0.070	12.5	0.012
RemeDi (++) Remask RL)	89.1	0.56	73.2	0.89	24.8	0.086

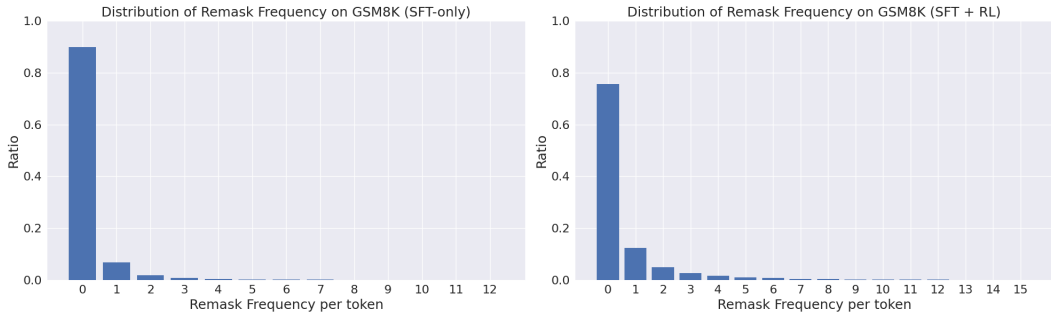


Figure 12: Comparison of per-token remask frequency between the SFT-only (left) and the RL-trained (right) model on GSM8K. The RL-trained model performs remasking more frequently on average, indicating that RL encourages more remasking behavior.

that the RL objective explicitly encourages the model to revise uncertain tokens more often, which correlates with the observed improvement in final answer quality.

Fig. 13 reports the remask ratios at different timesteps when inferencing on GSM8K. The ratio first increases and then steadily decreases as the denoising process converges. Since the number of unmasked tokens K_n is explicitly controlled at each step, the process naturally terminates without spikes of remasking in late stage.

Fig. 14 shows the throughput–performance trade-off of RemeDi compared to AR and DLM baselines. By adjusting the number of tokens denoised per step, RemeDi achieves acceleration in decoding speed with only small drops in GSM8K accuracy. It forms a better Pareto front than LLaDA, Dream, and other auto-regressive models.

Fig. 15 shows the reward curve for Remask-RL and LLaDOU-RL on GSM8K. Due to the larger action space introduced by the remask operation, Remask-RL starts with a lower initial reward. However, it quickly surpasses LLaDOU-RL within the early stages of training and maintains a consistently higher reward thereafter. This trend is consistent with the accuracy comparison in Table 5, where Remask-RL demonstrates both faster convergence and a higher final performance.

THE USE OF LARGE LANGUAGE MODELS

We used large language models (LLMs) in two limited ways: (1) for minor English polishing of the paper text, and (2) as required by the AlpacaEval benchmark, where LLMs are invoked for automatic evaluation.

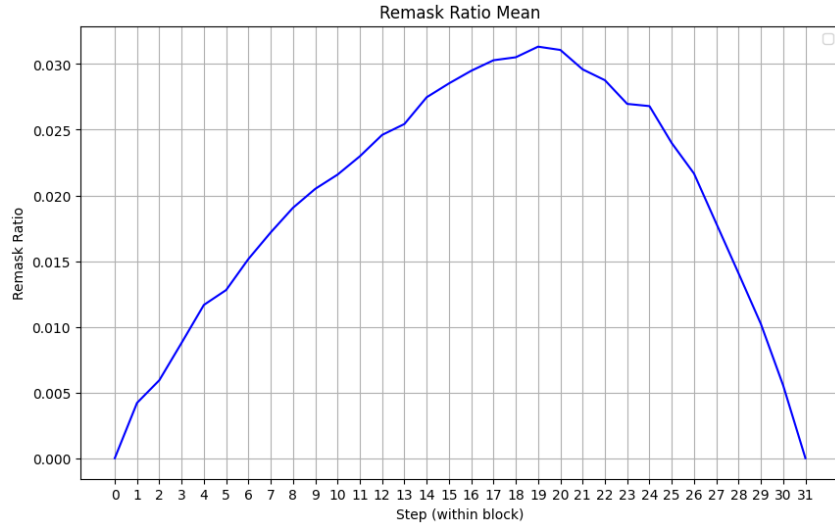


Figure 13: Average remask ratio across diffusion timesteps when inferencing on GSM8K. The remask ratio rises during early steps and then gradually declines as the process converges. Because the number of unmasked tokens K_n is explicitly defined at every diffusion step, the procedure ensures stable remask termination and avoids late-stage re-emergence of remasking.

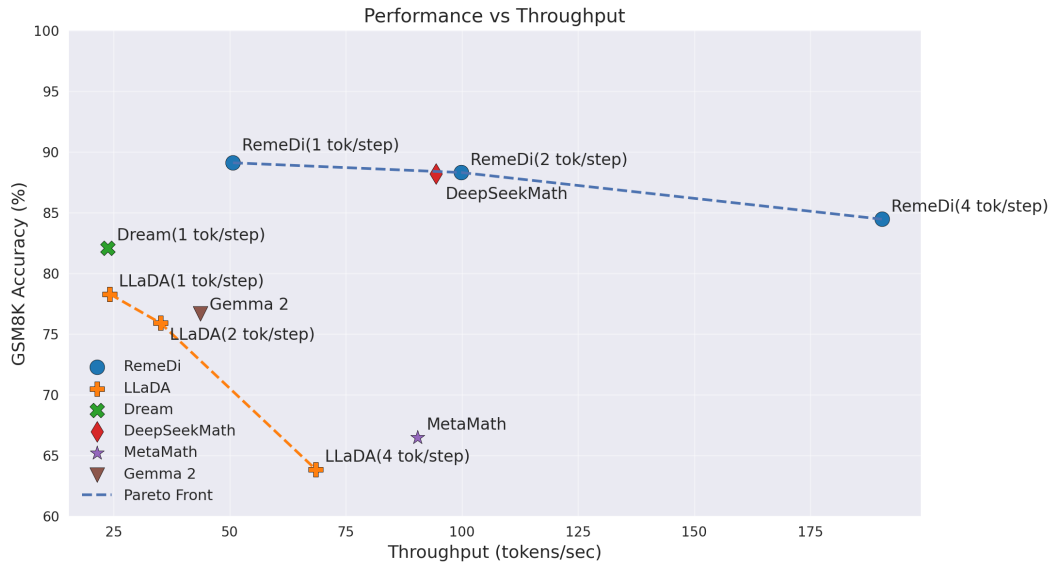


Figure 14: Throughput–performance trade-off of RemeDi compared with other AR and DLM models. By increasing the number of denoised tokens per step, RemeDi provides a smooth quality–latency trade-off. All results are measured with batch size 1 and sequence length 1024 on a single H800 GPU.

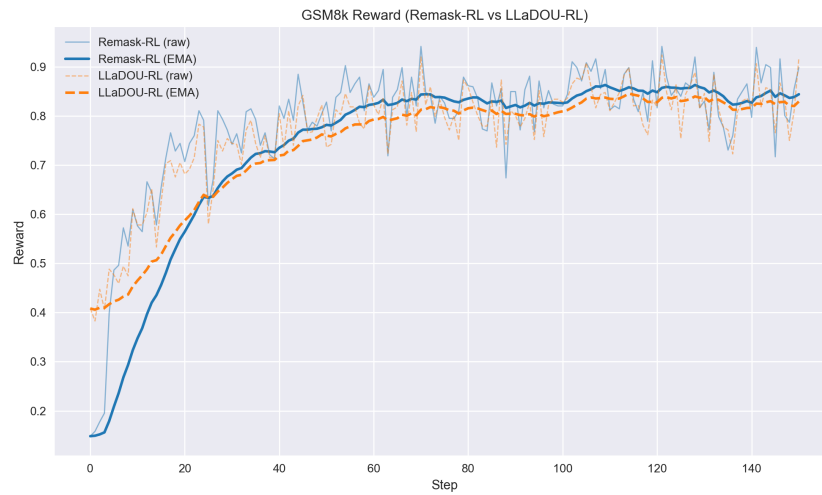


Figure 15: GSM8K training reward curves comparing Remask-RL and LLaDOU-RL. Solid lines show EMA-smoothed rewards, and faint lines denote raw step-wise values. Remask-RL exhibits faster early-stage improvement and ultimately reaches a higher reward, consistent with the performance gains summarized in Table 5