

# RFG: TEST-TIME SCALING FOR DIFFUSION LARGE LANGUAGE MODEL REASONING WITH REWARD-FREE GUIDANCE

Tianlang Chen \*    Minkai Xu \*    Jure Leskovec    Stefano Ermon  
Stanford University

## ABSTRACT

Diffusion Large Language Models (dLLMs) have shown great potential in language modeling, yet enhancing their capacity for complex reasoning remains a critical challenge. For autoregressive language models, this is typically addressed by guiding the reasoning process step by step using Process Reward Models (PRMs), which necessitate dense annotation for intermediate steps. However, this approach cannot be applied to dLLMs, since their intermediate generations are partially masked, non-sequential states rather than complete prefixes. Here we propose Reward-Free Guidance (RFG), a training-free framework that guides the reasoning trajectory of dLLMs without explicit process reward models. We provide theoretical justification that a process reward for partially masked states can be parameterized by the log-likelihood ratio of a policy and a reference model, which can be instantiated with off-the-shelf dLLM checkpoints without additional training. Extensive experiments demonstrate that RFG consistently outperforms various state-of-the-art post-trained dLLM baselines, achieving absolute accuracy gains of up to 9.2%.

## 1 INTRODUCTION

Diffusion Large Language Models (dLLMs) have recently achieved remarkable progress in language modeling (Sohl-Dickstein et al., 2015; Austin et al., 2021a; Campbell et al., 2022; Meng et al., 2022; Lou et al., 2024; Sahoo et al., 2024; Shi et al., 2024; Xu et al., 2025). By scaling up mask-predict pretraining via bidirectional modeling, dLLMs have demonstrated competitive capabilities across diverse domains, from mathematical reasoning and planning (Gong et al., 2023; Zhao et al., 2025) to expert coding (Gong et al., 2025b). With unique advantages such as bidirectional context modeling and parallel decoding, dLLMs now rival or even surpass autoregressive (AR) baselines in specific settings (Prabhudesai et al., 2025). Despite this progress, current dLLM research primarily focus on resource-intensive pre-training or domain-specific post-training, with limited exploration of *test-time scaling*. Test-time scaling aims to design algorithms that effectively utilize additional inference compute to improve accuracy while circumventing expensive post-training (Lightman et al., 2023; Wang et al., 2023; Snell et al., 2024). In the AR LLM landscape, scaling test-time compute has proven essential for solving complex reasoning tasks (Brown et al., 2024; Muennighoff et al., 2025), enabling significant performance gain without the need for expensive supervised fine-tuning or reinforcement learning (RL). Consequently, developing effective methods to unlock similar test-time scaling capabilities for dLLMs remains a critical yet underexplored direction.

Standard approaches for test-time scaling in AR LLMs typically rely on an auxiliary reward model to guide the search process. While initial attempts utilized Outcome Reward Models (ORMs) for coarse-grained feedback on the whole generation (e.g., Best-of-N sampling (Cobbe et al., 2021; Lightman et al., 2023)), recent approaches have shifted toward fine-grained Process Reward Models (PRMs). PRMs assign scores to intermediate reasoning steps, enabling search algorithms to prune incorrect paths early (Wang et al., 2023; Lu et al., 2024). However, transferring this success to dLLMs is non-trivial. While the community has explored naive solutions analogous to the Best-of-N approach (Wang et al., 2025a; Singhal et al., 2025; Dang et al., 2025), it is challenging

---

\*Equal contribution.

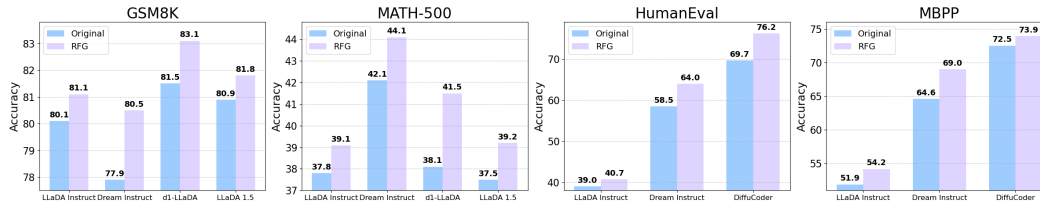


Figure 1: RFG consistently enhances performance over the original post-trained baseline across diverse benchmarks and model architectures. Results are averaged across different generation sequence lengths.

to directly adapt PRMs to dLLMs to guide the denoising process step by step. AR LLMs generates complete prefixes in a sequential manner which standard PRMs are designed to evaluate. In contrast, dLLM generation produces intermediate states comprised of partially masked sentences which are unintelligible to standard PRMs. Training a specialized PRM for these masked sentences would require collecting prohibitive amounts of dense, step-wise annotations on the masked data. This structural incompatibility highlights the need for a reward formulation native to the any-order generation nature of dLLMs.

Here we bridge this gap by proposing **Reward-Free Guidance (RFG)**, a principled framework for guiding dLLM reasoning without training reward models on dense annotations.<sup>1</sup> We view test-time scaling for dLLMs as a guided sampling problem, where guidance steers each denoising step toward accurate reasoning traces. Our key insight is that the reward of a complete trajectory can be parameterized as the log-likelihood ratio between an aligned policy model and a reference base model, effectively measuring how much the aligned model prefers the trajectory over the base model. Such reparameterization has been widely adopted in various RL and preference optimization literature. We theoretically demonstrate that this trajectory-level reward can be decomposed into an implicit step-wise process reward (c.f., Proposition 3.1), which takes the form of a log-likelihood ratio calculated on partially masked states. We thereby obtain an implicit process reward without training an explicit PRM. We then sample from a distribution reweighted by this implicit process reward, using a formulation (c.f., Equation 3) inspired by the well-studied Classifier-Free Guidance (CFG) (Ho & Salimans, 2022). Furthermore, we highlight that RFG is general and the implicit PRM can be instantiated using any off-the-shelf dLLM post-training methods without reliance on specific training objectives.

We validate RFG on challenging benchmarks spanning mathematical reasoning (GSM8K, MATH-500) and code generation (HumanEval, MBPP). RFG consistently yields significant gains over diverse post-trained baselines from state-of-the-art dLLM architectures including LLaDA (Nie et al., 2025b) and Dream (Ye et al., 2025), as shown in Figure 1. Crucially, RFG also outperforms compute-matched ensemble baselines, confirming that the improvements stem from our principled framework rather than merely increased compute. These results establish RFG as a robust, training-free framework for scaling test-time reasoning in non-autoregressive dLLMs.

## 2 PRELIMINARIES

**Discrete Diffusion Models.** Diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2020) are a class of generative models that learn a data distribution by reversing a progressive noising process. Originally developed for continuous data like images, the forward process gradually perturbs a clean sample into increasingly noisy states, eventually transforming it into pure noise. A neural network is then trained to reverse this process. Adapting this paradigm to discrete data such as text is non-trivial, as adding small amounts of Gaussian noise is not well-defined for discrete tokens. This led to the development of discrete diffusion models (Austin et al., 2021a), which define the forward process using a Markov chain that gradually randomizes discrete tokens, often converging to a uniform distribution. An effective and intuitive case of discrete diffusion is the

<sup>1</sup>We term our method "Reward-Free" following the convention of Classifier-Free Guidance (CFG). Just as CFG eliminates the need for an explicit classifier by deriving implicit gradients from the model itself, RFG eliminates the need for an explicit, separately trained reward model by deriving implicit rewards from the policy and reference models.

masked diffusion model (Campbell et al., 2022; Lou et al., 2024; Shi et al., 2024; Sahoo et al., 2024). In this framework, the "noise" is a special [MASK] token. The forward process progressively masks tokens in the input sequence until the entire sequence is masked. The reverse process then learns to predict the original tokens given a partially masked sequence. This formulation naturally leverages architectures and training objectives similar to masked language modeling (MLM) (Devlin et al., 2019), while inheriting the iterative refinement of diffusion.

**Diffusion Large Language Models.** Diffusion large language models (dLLMs) (Nie et al., 2025b; Ye et al., 2025; Inception Labs et al., 2025; DeepMind, 2025) scale masked diffusion models to large corpora and long reasoning sequences, serving as an alternative to autoregressive (AR) language models. The forward process of a dLLM takes the original sequence  $\mathbf{x}_0$  as input and progressively masks it following the distribution

$$q(\mathbf{x}_t|\mathbf{x}_0) = \prod_{i=0}^L q(\mathbf{x}_t^{(i)}|\mathbf{x}_0^{(i)}), \quad q(\mathbf{x}_t^{(i)}|\mathbf{x}_0^{(i)}) = \begin{cases} 1 - \alpha_t, & \mathbf{x}_t^{(i)} = [\text{MASK}] \\ \alpha_t, & \mathbf{x}_t^{(i)} = \mathbf{x}_0^{(i)} \end{cases}$$

where  $t \in [0, 1]$  is the timestep,  $L$  is the sequence length, and  $\alpha_t$  denotes a noise schedule that decreases monotonically with  $t$ , satisfying  $\alpha_0 = 1$  and  $\alpha_1 = 0$ . At timestep  $t$ , the ratio of masked token is  $1 - \alpha_t$ . A key property of dLLMs is that token prediction is not tied to a fixed left-to-right ordering. At each reverse step, the model can unmask any subset of masked positions conditioned on the visible context, enabling *any-order generation* in contrast to the strictly sequential nature of AR LLMs. Any-order generation reduces exposure bias, allows parallelized inference, and supports iterative refinement, making dLLMs well-suited for reasoning tasks that benefit from revisiting or correcting intermediate states. The detailed formulation of dLLMs is provided in Appendix B.

### 3 REWARD FREE GUIDANCE (RFG)

In this section, we elaborate on the details of our proposed RFG framework. We show how the reward for the denoising trajectory can be parameterized as the log-density ratio of two dLLMs, and then how an implicit PRM for each denoising step can be freely derived from the trajectory-level reward without additional training. Existing PRMs in AR LLM literature typically require fine-grained step labels, and then the PRM is trained to predict the quality of partially generated responses (Lightman et al., 2023; Wang et al., 2023; Lu et al., 2024). However, such labels are expensive to collect, and even impossible in the context of dLLMs since the intermediate generations are typically partial sentences with tokens masked in random positions.

To this end, we propose to *obtain a trajectory-level reward with reparameterization, and then freely derive an implicit PRM for each denoising step by decomposing the trajectory-level reward*. Formally, we have the following theoretical justification:

**Proposition 3.1.** *Given a diffusion trajectory-level reward that is parameterized as the log-likelihood ratio of two dLLMs, i.e.,  $r_\theta(\mathbf{x}_{0:T}) := \beta \log \frac{p_\theta(\mathbf{x}_{0:T})}{p_{\text{ref}}(\mathbf{x}_{0:T})}$ . Define  $Q_\theta^t(\mathbf{x}_{t-1}, \mathbf{x}_{t:T}) := \sum_{i=t}^T \beta \log \frac{p_\theta(\mathbf{x}_{i-1}|\mathbf{x}_{i:T})}{p_{\text{ref}}(\mathbf{x}_{i-1}|\mathbf{x}_{i:T})}$ , then we have that  $Q_\theta^t$  is the expectation of exponential  $r_\theta$  at step  $t$ :*

$$Q_\theta^t(\mathbf{x}_{t-1}, \mathbf{x}_{t:T}) = \beta \log \mathbb{E}_{p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{x}_{t-1:T})} e^{\frac{1}{\beta} r_\theta(\mathbf{x}_{0:T})}, \quad (1)$$

where  $\beta$  is a hyperparameter for weighting the reward function.

The full proof of the proposition is provided in Appendix A. The theoretical justification indicates that when parameterizing  $r_\theta(\mathbf{x}_{0:T})$  as the log-likelihood ratio,  $Q_\theta^t$  can be viewed as Q function representing the expectation of the overall trajectory reward  $r_\theta(\mathbf{x}_{0:T})$  at step  $t$ . Following the classic setup in RL literature, the step reward  $r_\theta^t$  can be written as:

$$r_\theta^t(\mathbf{x}_{t-1}|\mathbf{x}_t) = Q_\theta^t - Q_\theta^{t+1} = \beta \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_{t:T})}{p_{\text{ref}}(\mathbf{x}_{t-1}|\mathbf{x}_{t:T})} = \beta \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{p_{\text{ref}}(\mathbf{x}_{t-1}|\mathbf{x}_t)}. \quad (2)$$

Notably, with such PRMs, we can then conduct reward-guided sampling from dLLMs by reweighted denoising transitions, where the log probability can be written as:

$$\begin{aligned} \log p^*(\mathbf{x}_{t-1}|\mathbf{x}_t) &= \log p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) e^{\frac{1}{\beta} r_\theta^t(\mathbf{x}_{t-1}|\mathbf{x}_t)} + C \\ &= (1+w) \log p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) - w \log p_{\text{ref}}(\mathbf{x}_{t-1}|\mathbf{x}_t) + C, \end{aligned} \quad (3)$$

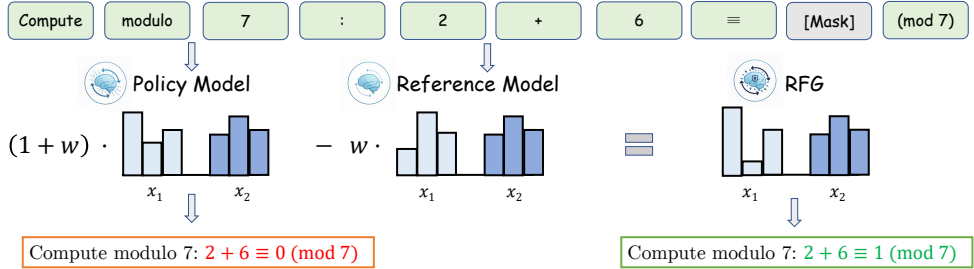


Figure 2: Illustration of RFG sampling. By linearly combining logits from the policy and reference models at each denoising step, RFG steers the generation trajectory toward accurate reasoning traces, resulting in improved performance without external reward models.

where  $w = \frac{\beta}{\gamma}$  is a new hyperparameter to control the guidance strength and  $C$  is a normalizing constant. Trivially, setting  $w = -1$  or  $w = 0$  can recover the exact sampling of original  $p_{\text{ref}}$  or  $p_{\theta}$ , respectively, with  $w > 0$  acting an over-emphasizing behavior of  $p_{\theta}$ . Such formulation provides a vital insight: by simply parameterizing a reward model on the whole diffusion sampling trajectory, we can freely obtain PRMs for each denoising step without any additional training. In the later paragraphs, we will discuss the connections between our RFG framework and other diffusion sampling guidance methods, and provide explanations on how we can obtain the trajectory reward model using any off-the-shelf RL or instruction-tuned checkpoints.

**Connections to diffusion guidance method.** A widely adopted technique for guiding diffusion sampling is *classifier-free guidance* (CFG) (Ho & Salimans, 2022), which pushes samples towards high class-confidence regions by reweighted denoising steps. In this section, we revisit our proposed RFG from the CFG perspective, and provide a holistic view of their connections.

Generally, CFG guidance in diffusion models involves two diffusion models  $p_{\text{unconditional}}(\mathbf{x}; \mathbf{c})$  and  $p_{\text{conditional}}(\mathbf{x})$ . In the original CFG in Gaussian diffusion, the guidance is achieved by *extrapolating* between the scores of two denoising models with a weight  $w$ :

$$\nabla_{\mathbf{x}} \log p^*(\mathbf{x}|\mathbf{c}) = \nabla_{\mathbf{x}} \log p_{\text{conditional}}(\mathbf{x}|\mathbf{c}) + w \nabla_{\mathbf{x}} \log \frac{p_{\text{conditional}}(\mathbf{x}|\mathbf{c})}{p_{\text{unconditional}}(\mathbf{x})}.$$

The effect of such guidance is standard sampling from  $p_{\text{conditional}}$  plus a drifting term to shift the sampling direction (for  $w > 0$ ) towards the ratio of  $p_{\text{conditional}}$  and  $p_{\text{unconditional}}$ . This ratio comes from the Bayesian rule that:

$$p^*(\mathbf{x}|\mathbf{c}) \propto p_{\text{conditional}}(\mathbf{x}|\mathbf{c}) \cdot p(\mathbf{c}|\mathbf{x})^w, \quad \text{where } p(\mathbf{c}|\mathbf{x}) \propto \frac{p_{\text{conditional}}(\mathbf{x}|\mathbf{c})}{p_{\text{unconditional}}(\mathbf{x})}.$$

The key idea of CFG is to reparameterize a hypothetical classifier as the likelihood ratio of conditional and unconditional diffusion models, which encourages the model to draw samples from density  $p_{\text{conditional}}$  over  $p_{\text{unconditional}}$ . From this perspective, there comes a clear connection between RFG and CFG: CFG guides the sampling with classifier  $p(\mathbf{c}|\mathbf{x})$  reparameterized as  $\frac{p_{\text{conditional}}(\mathbf{x}|\mathbf{c})}{p_{\text{unconditional}}(\mathbf{x})}$ , while RFG reweights the sampling with the reward  $r(\mathbf{x})$  that is reparameterized as  $\frac{p_{\theta}(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})}$ .

**Implementation of  $r_{\theta}$ .** We emphasize that Proposition 3.1 is general and agnostic to any training method of the reward model. Specifically, by off-the-shelf dLLMs fine-tuned with RL and preference optimization methods, we have the conclusion that the optimal policy will converge to  $p_{\theta}(\mathbf{x}) = p_{\text{ref}}(\mathbf{x})e^{\frac{1}{\beta}r(\mathbf{x})}$  after the training (Peters & Schaal, 2007). Directly rearranging the converged policy, we have that  $r_{\theta}(\mathbf{x}) = \beta \log \frac{p_{\theta}(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})}$ , which has already met our needs in guided sampling. Therefore, in practice, we can take any pair of a pretrained dLLM as the reference  $p_{\text{ref}}$  and a post-trained one via RL on certain tasks as the policy  $p_{\theta}$ , and conduct guided sampling via RFG framework. Interestingly, from an empirical perspective, we also observe that even without explicit RL, taking an instruction fine-tuned model as  $p_{\theta}$  can also offer significant performance gain, showing the generalization of RFG beyond our theoretical form.

**Algorithm 1** Sampling with Reward-Free Guidance (RFG)

**Require:** Reference model  $p_{\text{ref}}$ , policy model  $p_{\theta}$ , guidance strength  $w$ , query  $q$ , answer length  $L$ , sampling steps  $N$ , denoising strategy  $\mathcal{S}$

- 1: Initialize  $\mathbf{x}_N \leftarrow [\text{MASK}]^L$  ▷ Start with a fully masked sequence of length  $L$
- 2: **for**  $t \leftarrow N$  **down to** 1 **do**
- 3:    $\log \pi_{\text{ref}} \leftarrow \text{Logits}(p_{\text{ref}}(\cdot | \mathbf{x}_t, q))$  ▷ Get logits from reference model
- 4:    $\log \pi_{\theta} \leftarrow \text{Logits}(p_{\theta}(\cdot | \mathbf{x}_t, q))$  ▷ Get logits from enhanced model
- 5:    $\log \pi_{\text{RFG}} \leftarrow (1 + w) \log \pi_{\theta} - w \log \pi_{\text{ref}}$  ▷ Combine logits via RFG
- 6:    $\mathbf{x}_{t-1} \leftarrow \mathcal{S}(\mathbf{x}_t, \log \pi_{\text{RFG}}, t/N)$  ▷ Generate the next state using the guided logits
- 7: **end for**
- 8: **return**  $\mathbf{x}_0$

Table 1: Performance on GSM8K across different sequence lengths. Best in **bold**. Gain over the original post-trained model in **green**.

Model	Seq Len = 128				Seq Len = 256				Seq Len = 512			
	Original	Ensemble	RFG	Gain	Original	Ensemble	RFG	Gain	Original	Ensemble	RFG	Gain
<b>Instruction Fine-tuning</b>												
LLaDA 8B Instruct	76.3	39.3	<b>77.2</b>	<b>+0.9</b>	79.8	68.0	<b>81.3</b>	<b>+1.5</b>	84.1	80.1	<b>84.7</b>	<b>+0.6</b>
Dream 7B Instruct	67.9	27.7	<b>74.2</b>	<b>+6.3</b>	80.9	68.3	<b>82.1</b>	<b>+1.2</b>	84.8	84.0	<b>85.1</b>	<b>+0.3</b>
<b>Reinforcement Learning</b>												
d1-LLaDA	79.3	42.4	<b>79.8</b>	<b>+0.5</b>	82.5	68.5	<b>84.7</b>	<b>+2.2</b>	82.8	80.1	<b>84.7</b>	<b>+1.9</b>
LLaDA 1.5	76.3	45.3	<b>78.0</b>	<b>+1.7</b>	81.6	70.9	<b>82.1</b>	<b>+0.5</b>	84.7	80.5	<b>85.4</b>	<b>+0.7</b>

Table 2: Performance on MATH-500 across different sequence lengths. Best in **bold**. Gain over the original post-trained model in **green**.

Model	Seq Len = 128				Seq Len = 256				Seq Len = 512			
	Original	Ensemble	RFG	Gain	Original	Ensemble	RFG	Gain	Original	Ensemble	RFG	Gain
<b>Instruction Fine tuning</b>												
LLaDA 8B Instruct	32.4	22.8	<b>34.2</b>	<b>+1.8</b>	38.4	33.6	<b>39.6</b>	<b>+1.2</b>	42.6	37.6	<b>43.6</b>	<b>+1.0</b>
Dream 7B Instruct	32.2	11.6	<b>33.2</b>	<b>+1.0</b>	43.6	30.6	<b>46.4</b>	<b>+2.8</b>	50.4	51.2	<b>52.6</b>	<b>+2.2</b>
<b>Reinforcement Learning</b>												
d1 LLaDA	34.0	25.2	<b>37.4</b>	<b>+3.4</b>	38.8	34.8	<b>41.6</b>	<b>+2.8</b>	41.6	38.8	<b>45.4</b>	<b>+3.8</b>
LLaDA 1.5	32.0	23.0	<b>34.0</b>	<b>+2.0</b>	38.6	34.2	<b>39.6</b>	<b>+1.0</b>	42.0	40.2	<b>44.0</b>	<b>+2.0</b>

Table 3: Performance on HumanEval across different sequence lengths. Best in **bold**. Gain over the original post-trained model in **green**.

Model	Seq Len = 128				Seq Len = 256				Seq Len = 512			
	Original	Ensemble	RFG	Gain	Original	Ensemble	RFG	Gain	Original	Ensemble	RFG	Gain
<b>Instruction Fine tuning</b>												
LLaDA 8B Instruct	32.3	18.3	<b>33.5</b>	<b>+1.2</b>	39.6	37.2	<b>40.9</b>	<b>+1.3</b>	45.1	40.9	<b>47.6</b>	<b>+2.5</b>
Dream 7B Instruct	52.4	48.2	<b>60.4</b>	<b>+8.0</b>	61.0	59.1	<b>66.5</b>	<b>+5.5</b>	62.2	58.5	<b>65.2</b>	<b>+3.0</b>
<b>Reinforcement Learning</b>												
diffuCoder	70.1	64.6	<b>73.8</b>	<b>+3.7</b>	69.5	72.6	<b>76.2</b>	<b>+6.7</b>	69.5	73.2	<b>78.7</b>	<b>+9.2</b>

Table 4: Performance on MBPP across different sequence lengths. Best in **bold**. Gain over the original post-trained model in **green**.

Model	Seq Len = 128				Seq Len = 256				Seq Len = 512			
	Original	Ensemble	RFG	Gain	Original	Ensemble	RFG	Gain	Original	Ensemble	RFG	Gain
<b>Instruction Fine tuning</b>												
LLaDA 8B Instruct	51.0	52.5	<b>54.1</b>	<b>+3.1</b>	54.9	51.4	<b>55.6</b>	<b>+0.7</b>	49.8	49.8	<b>52.9</b>	<b>+3.1</b>
Dream 7B Instruct	65.0	65.8	<b>68.1</b>	<b>+3.1</b>	64.2	67.3	<b>68.5</b>	<b>+4.3</b>	64.6	69.6	<b>70.4</b>	<b>+5.8</b>
<b>Reinforcement Learning</b>												
diffuCoder	72.4	70.0	<b>73.9</b>	<b>+1.5</b>	72.4	72.8	<b>73.5</b>	<b>+1.1</b>	72.8	70.8	<b>74.3</b>	<b>+1.5</b>

**Sampling with RFG.** The practical implementation of RFG sampling involves a straightforward modification of the standard dLLM sampling loop at test time (Figure 2). RFG combines the logits from the policy and the reference model with guidance strength  $w$  towards the guided logit distribution, which is then used to determine the next denoised state  $\mathbf{x}_{t-1}$  (Algorithm 1). The denoising strategy  $\mathcal{S}$  is a generic function representing two key schedules in the dLLM denoising process: (1) the unmask schedule that decides *which* and *how many* [MASK] tokens in  $\mathbf{x}_t$  to recover given timestep  $t$  and the logits for all masked positions; (2) selecting the new tokens based on a chosen decoding strategy (*e.g.*, nucleus sampling) applied to the logits at the chosen position. In practice,  $p_{\theta}$  can be instantiated by any off-the-shelf post-trained model without any additional training.

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETUP

**Tasks.** We evaluate our method on four challenging reasoning benchmarks spanning **mathematical reasoning** and **code generation**. For mathematical reasoning, we use GSM8K (Cobbe et al., 2021), a benchmark of grade-school arithmetic word problems that requires multi-step symbolic reasoning, and MATH-500 (Hendrycks et al., 2021; Lightman et al., 2023), a curated set of challenging competition-level mathematics problems. For code generation, we use HumanEval (Chen et al., 2021), which contains handwritten Python programming problems described in docstrings, and the MBPP (Austin et al., 2021b) (sanitized), consisting of everyday Python tasks with natural language prompts and associated unit tests.

**Models.** We leverage two state-of-the-art families of dLLMs: LLaDA (Nie et al., 2025b) and Dream (Ye et al., 2025). We use LLaDA-Base and Dream-Base as the reference model, respectively. For each family, we use two categories of post-trained models as the policy model: instruction fine-tuned and RL-enhanced. For instruction fine-tuned models, we employ **LLaDA-Instruct** and **Dream-Instruct**. For RL-enhanced models, we include **d1-LLaDA** (Zhao et al., 2025), which applies GRPO (Shao et al., 2024) to enhance LLaDA on mathematical and logical tasks; **LLaDA-1.5** (Zhu et al., 2025), which introduces VRPO techniques to reduce variance when applying DPO (Rafailov et al., 2023) to LLaDA; and **diffuCoder** (Gong et al., 2025b), which builds on the Dream backbone and leverages coupled GRPO for code generation tasks.

**Baselines.** We compare RFG against two sets of baselines to rigorously evaluate its effectiveness. (1) **Original Post-Trained Model**, which directly apply the original post-trained model for inference. (2) **Naive Ensemble**, which computes the final logits by taking the average of logits of both post-trained (policy) and base (reference) models at each step. This baseline operates under an *identical compute budget* to RFG, serving as a controlled ablation. *Note:* We don’t include baselines using ORMs like Best-of-N (BoN) as a compute-matched comparison ( $N = 2$ ) lacks sufficient coverage. Furthermore, state-of-the-art dLLMs require near-zero sampling temperatures for optimal performance, which collapses generation diversity. Under these conditions, BoN yields results effectively identical to the original post-trained model.

**Evaluation.** For all benchmarks, we evaluate models under a zero-shot setting to assess their intrinsic reasoning capabilities without task-specific examples. We report accuracy for math reasoning tasks and pass@1 for code generation tasks. All results are reported on the official test sets of each benchmark. We use official checkpoints for all models whenever publicly available; for d1-LLaDA, whose checkpoint has not yet been released, we reproduce the model with the official source code. To ensure a fair comparison, all baselines are implemented and evaluated under the identical inference setting with the same hyperparameters. Additional implementation details and hyperparameters are provided in Appendix C.

### 4.2 MAIN RESULTS

We evaluate RFG on four benchmarks across different generation sequence length. Results are presented in Table 1, Table 2, Table 3, and Table 4. For every post-trained model variant, RFG consistently outperforms the original model across all tasks and sequence lengths. This demonstrates that RFG is highly effective at guiding the generation process towards more accurate and logically sound reasoning traces. Moreover, RFG yields significant gains over the naive ensemble baseline, which operates under an identical compute budget. This confirms that the performance improvements comes from our principled formulation rather than merely from incorporating the reference model’s parameters or more compute.

### 4.3 ANALYSIS AND DISCUSSION

**Qualitative Analysis.** We present qualitative examples from all models on every benchmark to illustrate how RFG improves reasoning quality in Appendix D due to space limitation. On mathematical reasoning tasks, RFG produces more coherent multi-step derivations and avoids halluci-

Table 5: Performance on mathematical reasoning tasks with a sequence length of 256. This setting doubles the number of tokens generated per step compared to Table 1 and Table 2.

Model	GSM8K				MATH-500			
	Original	Ensemble	RFG	Gain	Original	Ensemble	RFG	Gain
<b>Instruction Fine-tuning</b>								
LLaDA 8B Instruct	78.6	64.3	<b>80.7</b>	<b>+2.1</b>	34.2	29.4	<b>38.0</b>	<b>+3.8</b>
Dream 7B Instruct	68.5	58.6	<b>73.6</b>	<b>+5.1</b>	35.4	27.4	<b>37.6</b>	<b>+2.2</b>
<b>Reinforcement Learning</b>								
d1 LLaDA	79.8	67.3	<b>82.0</b>	<b>+2.2</b>	36.6	29.4	<b>39.2</b>	<b>+2.6</b>
LLaDA 1.5	81.5	70.2	<b>82.1</b>	<b>+0.6</b>	34.2	29.8	<b>38.0</b>	<b>+3.8</b>

Table 6: Performance on code generation tasks with a sequence length of 512. This setting doubles the number of tokens generated per step compared to Table 3 and Table 4.

Method	HumanEval				MBPP			
	Original	Ensemble	RFG	Gain	Original	Ensemble	RFG	Gain
<b>Instruction Fine-tuning</b>								
LLaDA 8B Instruct	37.2	32.9	<b>38.4</b>	<b>+1.2</b>	43.6	39.7	<b>44.4</b>	<b>+0.8</b>
Dream 7B Instruct	37.8	39.0	<b>42.7</b>	<b>+4.9</b>	48.2	52.1	<b>55.3</b>	<b>+7.1</b>
<b>Reinforcement Learning</b>								
diffuCoder	57.3	51.2	<b>60.4</b>	<b>+3.1</b>	61.1	36.6	<b>63.0</b>	<b>+1.9</b>

nation or inconsistent conclusions. On code generation tasks, RFG generates code that is not only syntactically correct but also more robust, reducing common errors such as missing edge conditions or incomplete logic. This qualitative case study corroborates our quantitative findings, illustrating how RFG effectively corrects and refines the reasoning trajectory during inference.

**Robustness to Accelerated Generation.** We further investigate the robustness of RFG in an accelerated generation setting, where models generate multiple tokens per step (i.e., fewer total denoising steps) to trade accuracy for inference speed. Specifically, we double the number of tokens generated per step (effectively halving the total denoising steps) compared to previous experiments, with results presented in Table 5 and Table 6. This serves as a challenging ablation of RFG’s stability. We find that RFG maintains substantial improvements over the baselines even in these coarser generation regimes, highlighting its practical effectiveness for latency-sensitive applications.

**Sensitivity to Guidance Strength.** We conduct a sensitivity analysis by varying the guidance strength hyperparameter  $w$ , with result presented in Figure 3. This demonstrate the robustness of RFG, as the improvement is not confined to a narrow, fine-tuned peak; instead, the model exhibits a wide plateau of strong performance across a broad range of  $w$  values. This indicates that RFG can achieve substantial gains without intensive hyperparameter tuning.

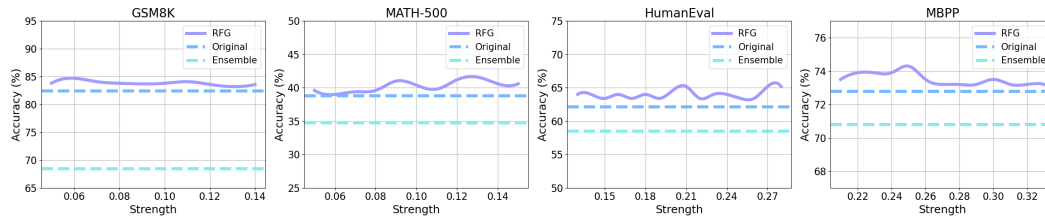


Figure 3: Accuracy under varying guidance strength  $w$ . GSM8K and MATH-500 uses d1-LLaDA, HumanEval and MBPP uses Dream-Instruct and DiffuCoder, respectively. RFG consistently improves performance over a broad range of guidance strength.

**Interpreting the Guidance Direction and Strength.** The RFG formulation can be written in the equivalent form:  $\log \pi_{\text{RFG}} = \log \pi_{\theta} + w(\log \pi_{\theta} - \log \pi_{\text{ref}})$ . With this form, RFG can be interpreted as steering the post-trained model’s distribution,  $\log \pi_{\theta}$ , along an optimization direction defined by the difference  $(\log \pi_{\theta} - \log \pi_{\text{ref}})$  with strength  $w$ . While RFG is motivated by theoretical derivation on RL models, our empirical results show RFG works well for SFT models too, and this interpretation helps explain why.

## 5 RELATED WORK

**Diffusion Language Models.** Early attempts at text diffusion relied on operating in continuous spaces, either by relaxing text tokens in a continuous form (Li et al., 2022; Dieleman et al., 2022; Gong et al., 2023; Chen et al., 2023; Wu et al., 2023) or by diffusing the continuous parameters of discrete distributions (Lou & Ermon, 2023; Lin et al., 2023; Graves et al., 2023; Xue et al., 2024). Despite conceptual simplicity, these approaches faced challenges in scalability (Gulrajani & Hashimoto, 2023). Alternatively, Austin et al. (2021a) introduced discrete diffusion that operates directly on discrete text tokens, leading to a proliferation of variants (Hoogeboom et al., 2021; He et al., 2022; Campbell et al., 2022; Meng et al., 2022; Sun et al., 2023; Gat et al., 2024; Sahoo et al., 2024; Shi et al., 2024; Zheng et al., 2025; Ou et al., 2025; Nie et al., 2025a; Gong et al., 2025a). In terms of generation order, Arriola et al. (2025) proposed Block Diffusion which generates block-by-block autoregressively while applying parallel diffusion within each block. Building upon these foundations, a significant breakthrough was the successful scaling of discrete diffusion language models, with LLaDA (Nie et al., 2025b) and Dream (Ye et al., 2025) demonstrating performance comparable to their autoregressive counterparts. To further boost reasoning and alignment, reinforcement learning has been applied. For instance, Zhao et al. (2025), Yang et al. (2025), and Tang et al. (2025) adapted the GRPO (Shao et al., 2024) objectives for dLLMs, while Zhu et al. (2025) introduced unbiased variance reduction techniques when applying DPO (Rafailov et al., 2023) to dLLMs. The scope has also expanded beyond text to multimodal domains, with models such as LaViDa (Li et al., 2025), MMaDA (Yang et al., 2025), and Dimple (Yu et al., 2025) integrating text diffusion with vision capabilities. SDAR (Cheng et al., 2025) introduces a hybrid paradigm that turns a pretrained autoregressive model into a blockwise diffusion model through a lightweight adaptation phase. TraDo (Wang et al., 2025b) presents TraceRL, a trajectory aware reinforcement learning framework that aligns the post training objective with the model’s inference trajectory.

**Steering Generative Models Reasoning.** Steering large generative models toward desired behaviors is a fundamental problem. In diffusion, guidance methods steer the generative process by adjusting the score or noise prediction with auxiliary signals such as class labels, weaker models, or reward gradients to bias sampling toward desired outputs without retraining. Classifier guidance (Dhariwal & Nichol, 2021) uses gradients from an external classifier to steer the diffusion process toward a target label, while classifier-free guidance (CFG) (Ho & Salimans, 2022) combines conditional and unconditional diffusion predictions to achieve the same effect without requiring a separate classifier. Karras et al. (2024) extended CFG and proposed autoguidance that guides the generation process with a deliberately less-trained and smaller version of itself rather than an unconditional model, decoupling quality from diversity and yielding state-of-the-art image generation. In contrast, RFG is entirely training-free: it does not deliberately train a weaker model and instead reuses off-the-shelf dLLMs that directly steer the reverse process. Similar ideas to autoguidance have also been explored in autoregressive LLMs, such as contrastive decoding (Li et al., 2023), which mitigates repetitiveness in generation. Nisonoff et al. (2025) applied guidance to discrete diffusion models, focusing on applications to scientific data such as molecules, DNA, and proteins. Another line of work uses reward-based guidance that typically trains an external reward model to score outputs, which has been successful in aligning autoregressive LLMs (Ouyang et al., 2022; Rafailov et al., 2023; Shao et al., 2024). Outcome reward models (ORMs) evaluate final answers without shaping intermediate reasoning, whereas process reward models (PRMs) assign rewards to step-level traces (Uesato et al., 2022; Lightman et al., 2023; Wang et al., 2023; Lu et al., 2024).

## 6 CONCLUSION

We introduced Reward-Free Guidance (RFG), a novel framework for enhancing the reasoning capabilities of diffusion large language models at test time. RFG guides the denoising process without requiring an explicitly trained process reward model. By parameterizing the reward as the log-likelihood ratio of a policy and reference dLLM, RFG elegantly connects a trajectory-level reward with step-wise guidance. Our theoretical analysis shows that RFG’s sampling process is equivalent to reward-guided sampling. Comprehensive experiments empirically validate the effectiveness of RFG. We envision this framework as a foundation for broader alignment and reasoning improvements in generative models, including multimodal diffusion and agentic reasoning systems, where test-time guidance offers a scalable and general alternative to costly retraining.

## REFERENCES

- Marianne Arriola, Aaron Gokaslan, Justin T Chiu, Zhihan Yang, Zhixuan Qi, Jiaqi Han, Subham Sekhar Sahoo, and Volodymyr Kuleshov. Block diffusion: Interpolating between autoregressive and diffusion language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://arxiv.org/abs/2503.09573>.
- Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in neural information processing systems*, 34:17981–17993, 2021a.
- 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*, 2021b.
- Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling. *arXiv preprint arXiv:2407.21787*, 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, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Ting Chen, Ruixiang ZHANG, and Geoffrey Hinton. Analog bits: Generating discrete data using diffusion models with self-conditioning. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=3itjR9QxFw>.
- Shuang Cheng, Yihan Bian, Dawei Liu, Yuhua Jiang, Yihao Liu, Linfeng Zhang, Wenhai Wang, Qipeng Guo, Kai Chen, Biqing Qi, et al. Sdar: A synergistic diffusion-autoregression paradigm for scalable sequence generation. *arXiv preprint arXiv:2510.06303*, 2025.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Meihua Dang, Jiaqi Han, Minkai Xu, Kai Xu, Akash Srivastava, and Stefano Ermon. Inference-time scaling of diffusion language models with particle gibbs sampling. *arXiv preprint arXiv:2507.08390*, 2025.
- DeepMind. Gemini diffusion. 2025. URL <https://deepmind.google/models/gemini-diffusion/>.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *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)*, June 2019.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- Sander Dieleman, Laurent Sartran, Arman Roshannai, Nikolay Savinov, Yaroslav Ganin, Pierre H Richemond, Arnaud Doucet, Robin Strudel, Chris Dyer, Conor Durkan, et al. Continuous diffusion for categorical data. *arXiv preprint arXiv:2211.15089*, 2022.
- Itai Gat, Tal Remez, Neta Shaul, Felix Kreuk, Ricky TQ Chen, Gabriel Synnaeve, Yossi Adi, and Yaron Lipman. Discrete flow matching. *Advances in Neural Information Processing Systems*, 37:133345–133385, 2024.

- Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, and Lingpeng Kong. Diffuseq: 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>.
- Shansan Gong, Shivam Agarwal, Yizhe Zhang, Jiacheng Ye, Lin Zheng, Mukai Li, Chenxin An, Peilin Zhao, Wei Bi, Jiawei Han, Hao Peng, and Lingpeng Kong. Scaling diffusion language models via adaptation from autoregressive models. In *The Thirteenth International Conference on Learning Representations*, 2025a. URL <https://openreview.net/forum?id=j1tSLYKwg8>.
- Shansan Gong, Ruixiang Zhang, Huangjie Zheng, Jiatao Gu, Navdeep Jaitly, Lingpeng Kong, and Yizhe Zhang. Diffucoder: Understanding and improving masked diffusion models for code generation. *arXiv preprint arXiv:2506.20639*, 2025b.
- Alex Graves, Rupesh Kumar Srivastava, Timothy Atkinson, and Faustino Gomez. Bayesian flow networks. *arXiv preprint arXiv:2308.07037*, 2023.
- Ishaan Gulrajani and Tatsunori B Hashimoto. Likelihood-based diffusion language models. *Advances in Neural Information Processing Systems*, 36:16693–16715, 2023.
- Zhengfu He, Tianxiang Sun, Kuanning Wang, Xuanjing Huang, and Xipeng Qiu. Diffusionbert: Improving generative masked language models with diffusion models. *arXiv preprint arXiv:2211.15029*, 2022.
- 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. *arXiv preprint arXiv:2103.03874*, 2021.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Emiel Hooeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows and multinomial diffusion: Learning categorical distributions. *Advances in neural information processing systems*, 34:12454–12465, 2021.
- Inception Labs, Samar Khanna, Siddhant Kharbanda, Shufan Li, Harshit Varma, Eric Wang, Sawyer Birnbaum, Ziyang Luo, Yanis Miraoui, Akash Palrecha, et al. Mercury: Ultra-fast language models based on diffusion. *arXiv preprint arXiv:2506.17298*, 2025.
- Tero Karras, Miika Aittala, Tuomas Kynkäänniemi, Jaakko Lehtinen, Timo Aila, and Samuli Laine. Guiding a diffusion model with a bad version of itself. *Advances in Neural Information Processing Systems*, 37:52996–53021, 2024.
- Shufan Li, Konstantinos Kallidromitis, Hritik Bansal, Akash Gokul, Yusuke Kato, Kazuki Kozuka, Jason Kuen, Zhe Lin, Kai-Wei Chang, and Aditya Grover. Lavidia: A large diffusion language model for multimodal understanding. *ArXiv preprint*, abs/2505.16839, 2025. URL <https://arxiv.org/abs/2505.16839>.
- Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-lm improves controllable text generation. *Advances in neural information processing systems*, 35:4328–4343, 2022.
- Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. Contrastive decoding: Open-ended text generation as optimization. In *The 61st Annual Meeting Of The Association For Computational Linguistics*, 2023.
- 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.

- Zhenghao Lin, Yeyun Gong, Yelong Shen, Tong Wu, Zhihao Fan, Chen Lin, Nan Duan, and Weizhu Chen. Text generation with diffusion language models: A pre-training approach with continuous paragraph denoise. In *International Conference on Machine Learning*, pp. 21051–21064. PMLR, 2023.
- Aaron Lou and Stefano Ermon. Reflected diffusion models. In *International Conference on Machine Learning*, pp. 22675–22701. PMLR, 2023.
- Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion modeling by estimating the ratios of the data distribution. In *International Conference on Machine Learning*, pp. 32819–32848. PMLR, 2024.
- Jianqiao Lu, Zhiyang Dou, Hongru Wang, Zeyu Cao, Jianbo Dai, Yunlong Feng, and Zhijiang Guo. Autopsv: Automated process-supervised verifier. *Advances in Neural Information Processing Systems*, 37:79935–79962, 2024.
- Chenlin Meng, Kristy Choi, Jiaming Song, and Stefano Ermon. Concrete score matching: Generalized score matching for discrete data. *Advances in Neural Information Processing Systems*, 35: 34532–34545, 2022.
- Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling, 2025. URL <https://arxiv.org/abs/2501.19393>.
- Shen Nie, Fengqi Zhu, Chao Du, Tianyu Pang, Qian Liu, Guangtao Zeng, Min Lin, and Chongxuan Li. Scaling up masked diffusion models on text. In *The Thirteenth International Conference on Learning Representations*, 2025a. URL <https://openreview.net/forum?id=WNvvwK0tut>.
- 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, 2025b. URL <https://arxiv.org/abs/2502.09992>.
- Hunter Nisonoff, Junhao Xiong, Stephan Allenspach, and Jennifer Listgarten. Unlocking guidance for discrete state-space diffusion and flow models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=XsgH154y07>.
- 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. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=sMyXP8Tanm>.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022.
- Jan Peters and Stefan Schaal. Reinforcement learning by reward-weighted regression for operational space control. In *Proceedings of the 24th international conference on Machine learning*, pp. 745–750, 2007.
- Mihir Prabhudesai, Mengning Wu, Amir Zadeh, Katerina Fragkiadaki, and Deepak Pathak. Diffusion beats autoregressive in data-constrained settings. *arXiv preprint arXiv:2507.15857*, 2025.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741, 2023.
- Subham Sahoo, Marianne Arriola, Yair Schiff, Aaron Gokaslan, Edgar Marroquin, Justin Chiu, Alexander Rush, and Volodymyr Kuleshov. Simple and effective masked diffusion language models. *Advances in Neural Information Processing Systems*, 37:130136–130184, 2024.

- 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.
- Jiaxin Shi, Kehang Han, Zhe Wang, Arnaud Doucet, and Michalis Titsias. Simplified and generalized masked diffusion for discrete data. *Advances in neural information processing systems*, 37: 103131–103167, 2024.
- Raghav Singhal, Zachary Horvitz, Ryan Teehan, Mengye Ren, Zhou Yu, Kathleen McKeown, and Rajesh Ranganath. A general framework for inference-time scaling and steering of diffusion models. *arXiv preprint arXiv:2501.06848*, 2025.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pp. 2256–2265. pmlr, 2015.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2020.
- Haoran Sun, Lijun Yu, Bo Dai, Dale Schuurmans, and Hanjun Dai. Score-based continuous-time discrete diffusion models. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=BYWWwSY2G5s>.
- 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.
- Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process-and outcome-based feedback. *arXiv preprint arXiv:2211.14275*, 2022.
- Guanghan Wang, Yair Schiff, Subham Sekhar Sahoo, and Volodymyr Kuleshov. Remasking discrete diffusion models with inference-time scaling. *arXiv preprint arXiv:2503.00307*, 2025a.
- Peiyi Wang, Lei Li, Zhihong Shao, RX Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. *arXiv preprint arXiv:2312.08935*, 2023.
- Yinjie Wang, Ling Yang, Bowen Li, Ye Tian, Ke Shen, and Mengdi Wang. Revolutionizing reinforcement learning framework for diffusion large language models. *arXiv preprint arXiv:2509.06949*, 2025b.
- Tong Wu, Zhihao Fan, Xiao Liu, Hai-Tao Zheng, Yeyun Gong, Jian Jiao, Juntao Li, Jian Guo, Nan Duan, Weizhu Chen, et al. Ar-diffusion: Auto-regressive diffusion model for text generation. *Advances in Neural Information Processing Systems*, 36:39957–39974, 2023.
- Minkai Xu, Tomas Geffner, Karsten Kreis, Weili Nie, Yilun Xu, Jure Leskovec, Stefano Ermon, and Arash Vahdat. Energy-based diffusion language models for text generation. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=sL2F9YCMXf>.
- Kaiwen Xue, Yuhao Zhou, Shen Nie, Xu Min, Xiaolu Zhang, Jun Zhou, and Chongxuan Li. Unifying bayesian flow networks and diffusion models through stochastic differential equations. In *International Conference on Machine Learning*, pp. 55656–55681. PMLR, 2024.
- Ling Yang, Ye Tian, Bowen Li, Xinchun Zhang, Ke Shen, Yunhai Tong, and Mengdi Wang. Mmada: Multimodal large diffusion language models. *ArXiv preprint*, abs/2505.15809, 2025. URL <https://arxiv.org/abs/2505.15809>.
- Jiacheng Ye, Zhihui Xie, Lin Zheng, Jiahui Gao, Zirui Wu, Xin Jiang, Zhenguo Li, and Lingpeng Kong. Dream 7b: Diffusion large language models, 2025.

- Runpeng Yu, Xinyin Ma, and Xinchao Wang. Dimple: Discrete diffusion multimodal large language model with parallel decoding. *ArXiv preprint*, abs/2505.16990, 2025. URL <https://arxiv.org/abs/2505.16990>.
- Lifan Yuan, Wendi Li, Huayu Chen, Ganqu Cui, Ning Ding, Kaiyan Zhang, Bowen Zhou, Zhiyuan Liu, and Hao Peng. Free process rewards without process labels. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=8ThnPFhGm8>.
- Siyan Zhao, Devaansh Gupta, Qinqing Zheng, and Aditya Grover. d1: Scaling reasoning in diffusion large language models via reinforcement learning. *ArXiv preprint*, abs/2504.12216, 2025. URL <https://arxiv.org/abs/2504.12216>.
- Kaiwen Zheng, Yongxin Chen, Hanzi Mao, Ming-Yu Liu, Jun Zhu, and Qinsheng Zhang. Masked diffusion models are secretly time-agnostic masked models and exploit inaccurate categorical sampling. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=CTC7CmirNr>.
- 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 PROOF

**Proposition 3.1.** *Given a diffusion trajectory-level reward that is parameterized as the log-likelihood ratio of two dLLMs, i.e.,  $r_\theta(\mathbf{x}_{0:T}) := \beta \log \frac{p_\theta(\mathbf{x}_{0:T})}{p_{\text{ref}}(\mathbf{x}_{0:T})}$ . Define  $Q_\theta^t(\mathbf{x}_{t-1}, \mathbf{x}_{t:T}) := \sum_{i=t}^T \beta \log \frac{p_\theta(\mathbf{x}_{i-1}|\mathbf{x}_{i:T})}{p_{\text{ref}}(\mathbf{x}_{i-1}|\mathbf{x}_{i:T})}$ , then we have that  $Q_\theta^t$  is the expectation of exponential  $r_\theta$  at step  $t$ :*

$$Q_\theta^t(\mathbf{x}_{t-1}, \mathbf{x}_{t:T}) = \beta \log \mathbb{E}_{p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{x}_{t-1:T})} e^{\frac{1}{\beta} r_\theta(\mathbf{x}_{0:T})}, \quad (4)$$

where  $\beta$  is a hyperparameter for weighting the reward function.

*Proof.* The proposition is mainly proven by induction. The proposition and proof are largely borrowed from related literature in autoregressive LLMs (Yuan et al., 2025). The key difference is that in previous work, the analyses are mainly for autoregressive generation, while in this paper, we focus on the sampling trajectory of dLLMs.

Suppose we are given dLLM that discrete the sampling trajectory into  $T$  steps. The proof of Theorem 3.1 can be decomposed into the two following arguments:

1. At  $t = 1$ ,  $Q_\theta^1(\mathbf{x}_0, \mathbf{x}_{1:T}) = r_\theta(\mathbf{x}_{0:T}) = \beta \log \mathbb{E}_{p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{x}_{0:T})} e^{\frac{1}{\beta} r_\theta(\mathbf{x}_{0:T})}$ ;
2. For  $\forall t \in \{1, \dots, T-1\}$ , if  $Q_\theta^t(\mathbf{x}_{t-1}, \mathbf{x}_{t:T}) = \beta \log \mathbb{E}_{p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{x}_{t-1:T})} e^{\frac{1}{\beta} r_\theta(\mathbf{x}_{0:T})}$ , then we would also have  $Q_\theta^{t+1}(\mathbf{x}_t, \mathbf{x}_{t+1:T}) = \beta \log \mathbb{E}_{p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{x}_t:T)} e^{\frac{1}{\beta} r_\theta(\mathbf{x}_{0:T})}$ .

**Proof of 1.** In dLLM, we have that  $p(\mathbf{x}_{0:T}) = p(x_T) \prod_{t=1}^T p(\mathbf{x}_{t-1}|\mathbf{x}_t)$ . Then we have

$$r_\theta(\mathbf{x}) := \beta \log \frac{p_\theta(\mathbf{x}_{0:T})}{p_{\text{ref}}(\mathbf{x}_{0:T})} = \beta \log \prod_{i=1}^T \frac{p_\theta(\mathbf{x}_{i-1}|\mathbf{x}_i)}{p_{\text{ref}}(\mathbf{x}_{i-1}|\mathbf{x}_i)} = \sum_{i=1}^T \beta \log \frac{p_\theta(\mathbf{x}_{i-1}|\mathbf{x}_i)}{p_{\text{ref}}(\mathbf{x}_{i-1}|\mathbf{x}_i)}.$$

Then we trivially have that:

$$Q_\theta^1(\mathbf{x}_0, \mathbf{x}_{1:T}) = \sum_{i=1}^T \beta \log \frac{p_\theta(\mathbf{x}_{i-1}|\mathbf{x}_i)}{p_{\text{ref}}(\mathbf{x}_{i-1}|\mathbf{x}_i)} = r_\theta(\mathbf{x}_{0:T}) = \beta \log \mathbb{E}_{p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{x}_{0:T})} e^{\frac{1}{\beta} r_\theta(\mathbf{x}_{0:T})}.$$

**Proof of 2.** For  $\forall t \in \{1, \dots, T-1\}$ , given  $Q_\theta^t(\mathbf{x}_{t-1}, \mathbf{x}_{t:T}) = \beta \log \mathbb{E}_{p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{x}_{t-1:T})} e^{\frac{1}{\beta} r_\theta(\mathbf{x}_{0:T})}$ , we have:

$$\begin{aligned} Q_\theta^{t+1}(\mathbf{x}_t, \mathbf{x}_{t+1:T}) &= \beta \sum_{i=t}^T \log \frac{p_\theta(\mathbf{x}_i|\mathbf{x}_{i+1:T})}{p_{\text{ref}}(\mathbf{x}_i|\mathbf{x}_{i+1:T})} \\ &= \beta \log \prod_{i=t}^T \frac{p_\theta(\mathbf{x}_i|\mathbf{x}_{i+1:T})}{p_{\text{ref}}(\mathbf{x}_i|\mathbf{x}_{i+1:T})} \\ &= \beta \log \prod_{i=t}^T \frac{p_\theta(\mathbf{x}_i|\mathbf{x}_{i+1:T})}{p_{\text{ref}}(\mathbf{x}_i|\mathbf{x}_{i+1:T})} \sum_{\mathbf{y}_{t-1}} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_{t:T}) \\ &= \beta \log \prod_{i=t}^T \frac{p_\theta(\mathbf{x}_i|\mathbf{x}_{i+1:T})}{p_{\text{ref}}(\mathbf{x}_i|\mathbf{x}_{i+1:T})} \sum_{\mathbf{y}_{t-1}} p_{\text{ref}}(\mathbf{x}_{t-1}|\mathbf{x}_{t:T}) \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_{t:T})}{p_{\text{ref}}(\mathbf{x}_{t-1}|\mathbf{x}_{t:T})} \\ &= \beta \log \prod_{i=t}^T \frac{p_\theta(\mathbf{x}_i|\mathbf{x}_{i+1:T})}{p_{\text{ref}}(\mathbf{x}_i|\mathbf{x}_{i+1:T})} \mathbb{E}_{p_{\text{ref}}(\mathbf{x}_{t-1}|\mathbf{x}_{t:T})} \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_{t:T})}{p_{\text{ref}}(\mathbf{x}_{t-1}|\mathbf{x}_{t:T})} \\ &= \beta \log \mathbb{E}_{p_{\text{ref}}(\mathbf{x}_{t-1}|\mathbf{x}_{t:T})} \prod_{i=t-1}^T \frac{p_\theta(\mathbf{x}_i|\mathbf{x}_{i+1:T})}{p_{\text{ref}}(\mathbf{x}_i|\mathbf{x}_{i+1:T})} \\ &= \beta \log \mathbb{E}_{p_{\text{ref}}(\mathbf{x}_{t-1}|\mathbf{x}_{t:T})} e^{\frac{1}{\beta} Q_\theta^t(\mathbf{x}_{t-1}, \mathbf{x}_{t:T})} \end{aligned}$$

$$\begin{aligned}
&= \beta \log \mathbb{E}_{p_{\text{ref}}(\mathbf{x}_{t-1}|\mathbf{x}_{t:T})} \mathbb{E}_{p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{x}_{t-1:T})} e^{\frac{1}{\beta} r_{\theta}(\mathbf{x}_{0:T})} \\
&= \beta \log \mathbb{E}_{p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{x}_{t:T})} e^{\frac{1}{\beta} r_{\theta}(\mathbf{x}_{0:T})}
\end{aligned}$$

which completes the proof.  $\square$

## B MASKED DIFFUSION LARGE LANGUAGE MODEL FORMULATION

We can frame a diffusion large language model (dLLM) as a scaled-up implementation of a masked diffusion model tailored for text generation. It operates on a sequence of tokens  $\mathbf{x}_0$  and learns to reverse a masking process.

**Forward Process.** The forward process  $q(\mathbf{x}_t|\mathbf{x}_{t-1})$  is a Markov process that replaces a subset of non-masked tokens in  $\mathbf{x}_{t-1}$  with a [MASK] token to produce  $\mathbf{x}_t$ . The number of tokens to mask at each step is determined by a predefined noise schedule  $\alpha_t$ , where  $t \in [0, 1]$  is the timestep. At timestep  $t$ , the ratio of masked token is  $1 - \alpha_t$ .  $\alpha_t$  strictly decreases with  $t$ , satisfying  $\alpha_0 = 1$  and  $\alpha_1 = 0$ . As an example, LLaDA (Nie et al., 2025b) adopts a linear schedule defined as  $\alpha_t = 1 - t$ . We then have the following the distribution

$$q_{t|0}(\mathbf{x}_t|\mathbf{x}_0) = \prod_{i=0}^L q_{t|0}(\mathbf{x}_t^{(i)}|\mathbf{x}_0^{(i)}), \quad q_{t|0}(\mathbf{x}_t^{(i)}|\mathbf{x}_0^{(i)}) = \begin{cases} 1 - \alpha_t, & \mathbf{x}_t^{(i)} = [\text{MASK}] \\ \alpha_t, & \mathbf{x}_t^{(i)} = \mathbf{x}_0^{(i)} \end{cases}$$

**Reverse Process.** The reverse process aims to predict the original tokens that were masked in  $\mathbf{x}_t$ . It is worth noting that during the forward process once a token is masked, it remains in the masked state and cannot transition to other states. Given this, the conditional distribution for the reverse process moving from a time step  $t$  to an earlier step  $s$ , where  $0 \leq s < t \leq 1$ , is expressed as

$$q_{s|t}(\mathbf{x}_s^{(i)}|\mathbf{x}_t) = \begin{cases} 1, & \mathbf{x}_t^{(i)} \neq [\text{MASK}], \mathbf{x}_s^{(i)} = \mathbf{x}_t^{(i)} \\ \frac{1 - \alpha_s}{1 - \alpha_t}, & \mathbf{x}_t^{(i)} = [\text{MASK}], \mathbf{x}_s^{(i)} = [\text{MASK}] \\ \frac{\alpha_s - \alpha_t}{1 - \alpha_t} q_{0|t}(\mathbf{x}_s^{(i)}|\mathbf{x}_t), & \mathbf{x}_t^{(i)} = [\text{MASK}], \mathbf{x}_s^{(i)} \neq [\text{MASK}] \\ 0, & \text{otherwise} \end{cases}$$

In practice,  $q_{0|t}(\mathbf{x}_s^{(i)}|\mathbf{x}_t)$  is approximated by the dLLM that recover the original token in sequence  $\mathbf{x}_0$  given the partially masked sequence  $\mathbf{x}_t$ .

**Training.** The dLLM, denoted as  $p_{\theta}(\cdot|\mathbf{x}_t)$ , is trained to reconstruct the original sequence  $\mathbf{x}_0$  by predicting all masked tokens in the sequence  $\mathbf{x}_t$ , analogous to the masked language modeling objective. The loss is an upper bound of the negative log-likelihood of the model distribution:

$$\mathcal{L}(\theta) := \mathbb{E}_{\mathbf{x}_0 \sim p_{\text{data}}(\mathbf{x}_0), \mathbf{x}_t \sim q_{t|0}(\cdot|\mathbf{x}_0)} \left[ \int_{t=0}^{t=1} \frac{\alpha'_t}{1 - \alpha_t} \sum_{i=1}^L \mathbb{I}[\mathbf{x}_t^{(i)} = [\text{MASK}]] \log p_{\theta}(\mathbf{x}_0^{(i)}|\mathbf{x}_t) \right],$$

where  $p_{\text{data}}$  is the distribution of training data. For models that employ linear noise schedule  $\alpha_t = 1 - t$  (e.g., LLaDA), the loss function simplifies to

$$\mathcal{L}(\theta) = -\mathbb{E}_{t \sim \mathcal{U}[0,1], \mathbf{x}_0 \sim p_{\text{data}}(\mathbf{x}_0), \mathbf{x}_t \sim q_{t|0}(\cdot|\mathbf{x}_0)} \left[ \frac{1}{t} \sum_{i=1}^L \mathbb{I}[\mathbf{x}_t^{(i)} = [\text{MASK}]] \log p_{\theta}(\mathbf{x}_0^i|\mathbf{x}_t) \right].$$

**Inference.** Inference starts with a sequence of  $L$  [MASK] tokens. The model then iteratively unmask the sequence over  $T$  steps. At each step  $t$ , the model  $p_{\theta}(\mathbf{x}_0|\mathbf{x}_t)$  predicts a full sequence. A key feature of dLLMs is their flexibility in the unmasking order. Instead of a fixed left-to-right generation, dLLMs can employ an *any-order* generation scheme. For instance, at each step, the model might only replace the [MASK] tokens for which it has the highest prediction confidence (i.e., the lowest entropy). This new, partially filled sequence becomes  $\mathbf{x}_{t-1}$ , and the process repeats. This allows the model to build the sequence from the easiest parts to the hardest, leveraging bidirectional context at every step.

## C IMPLEMENTATION DETAILS

**Hyperparameters.** The number of generation steps is set equal to the generation length for experiments in Table 1, 2, 3, and 4, and is set to half of the generation length for experiments in Table 5 and 6. For the LLaDA family (LLaDA-Instruct, d1-LLaDA, and LLaDA-1.5), we adopt block diffusion with a block length of 32 for all tasks. We apply low-confidence remasking for all LLaDA-based models. The temperature is set to 0 for LLaDA family, following the setting from the original paper. For the Dream family (Dream-Instruct and DiffuCoder), we use nucleus sampling with  $p = 0.95$  and a temperature of 0.1, following the recommendation from the original paper. To ensure fairness, we use exactly the same set of hyperparameters for the original models, the naive ensemble baselines, and our RFG.

**Task Prompts.** For GSM8k and MATH-500, we provide the original problem statement followed by an instruction asking the model to output the final answer enclosed in a `\boxed` command. For HumanEval and MBPP, we prepend a generation prefix: “Here is the code to solve this problem\n```\npython”. For MBPP specifically, we include only a single test case in the prompt, striking a balance between clarifying the task details and avoiding potential test-data leakage. We use exact same setting for all model types for each task.

## D QUALITATIVE EXAMPLES

## D.1 GSM8K

## D.1.1 LLaDA INSTRUCT

## Query (GSM8K)

Marissa is hiking a 12-mile trail. She took 1 hour to walk the first 4 miles, then another hour to walk the next two miles. If she wants her average speed to be 4 miles per hour, what speed (in miles per hour) does she need to walk the remaining distance?

## Original (LLaDA Instruct)

First, let's determine the total distance Marissa has already walked:

$$4 \text{ miles} + 2 \text{ miles} = 6 \text{ miles}$$

Next, let's calculate the total time she has spent walking:

$$1 \text{ hour} + 1 \text{ hour} = 2 \text{ hours}$$

Now, let's find the remaining distance she needs to walk:

$$12 \text{ miles} - 6 \text{ miles} = 6 \text{ miles}$$

Next, let's determine the remaining time she has:

$$4 \text{ hours} - 2 \text{ hours} = 2 \text{ hours}$$

Finally, let's calculate the speed she needs to walk the remaining distance:

$$\text{Speed} = \frac{6 \text{ miles}}{2 \text{ hours}} = 3 \text{ miles per hour}$$

Therefore, Marissa needs to walk the remaining distance at a speed of

3 miles per hour to achieve an average speed of 4 miles per hour.

## RFG (LLaDA Instruct)

To determine how much Christina will spend on gift bags, we need to follow these steps:

First, let's determine the total time Marissa has already spent hiking. She took 1 hour to walk the first 4 miles and another 1 hour to walk the next 2 miles. Therefore, the total time spent is:

$$1 + 1 = 2 \text{ hours}$$

Next, we know the total distance of the trail is 12 miles. She has already walked 4 miles + 2 miles = 6 miles. Therefore, the remaining distance is:

$$12 - 6 = 6 \text{ miles}$$

Marissa wants her average speed to be 4 miles per hour for the entire trail. Therefore, the total time she needs to complete the trail is:

$$\frac{12}{4} = 3 \text{ hours}$$

Since she has already spent 2 hours hiking, she has 1 hour left to complete the remaining 6 miles. Therefore, the speed she needs to walk the remaining distance is:

$$\frac{6}{1} = 6 \text{ miles per hour}$$

The final answer is:

6

## D.1.2 DREAM INSTRUCT

## Query (GSM8K)

Carla is downloading a 200 GB file. Normally she can download 2 GB/minute, but 40% of the way through the download, Windows forces a restart to install updates, which takes 20 minutes. Then Carla has to restart the download from the beginning. How long does it take to download the file?

## Original (Dream Instruct)

First, we calculate the time it takes to download 40% of the 200 GB file at a rate of 2 GB/minute. 40% of 200 GB is 80 GB, so it takes  $80 \text{ GB} / 2 \text{ GB/minute} = 40$  minutes to download 80 GB.

Next, we calculate the time it takes to download the remaining 60% of the file at the same rate. 60% of 200 GB is 120 GB, so it takes  $120 \text{ GB} / 2 \text{ GB/minute} = 60$  minutes to download 120 GB.

Finally, we add the time it takes to download the first 40% of the file, the time it takes to download the remaining 60%, and the time it takes to install updates. The total time is 40 minutes + 60 minutes + 20 minutes = 120 minutes.

Therefore, it takes Carla 120 minutes to download the file.

## RFG (Dream Instruct)

1. First, calculate the amount of data Carla has downloaded when the 40

$$200 \times 0.4 = 80 \text{ GB}$$

2. Since Carla downloads at a rate of 2 GB/minute, the time it takes to download 80 GB is:

$$\frac{80}{2} = 40 \text{ minutes}$$

3. After the restart, Carla has to download the entire 200 GB again. The time it takes to download 200 GB at a rate of 2 GB/minute is:

$$\frac{200}{2} = 100 \text{ minutes}$$

4. Add the time for the restart and the time to download the entire file again to find the total time:

$$40 + 20 + 100 = 160 \text{ minutes}$$

$$\boxed{160} \text{ minutes}$$

## D.1.3 D1-LLADA

## Query (GSM8K)

Candice put 80 post-it notes in her purse before she headed out to her job at the coffee shop. On her way, she stopped off at the store and purchased a package of Post-it notes; At work, she placed a single Post-it note on each of 220 different cups of coffee. If she had 23 post-it notes remaining overall, how many Post-it notes were in the package that she purchased?

## Original (d1 LLaDA)

To determine how many Post-it notes were in the package that Candice purchased, we need to follow these steps:

1. Calculate the total number of Post-it notes Candice used at work. 2. Add the number of Post-it notes she had remaining to the number of Post-it notes she initially had in her purse. 3. Subtract the number of Post-it notes she used at work from the total number of Post-it notes she had after purchasing the package.

First, we know that Candice had 80 Post-it notes in her purse initially and 23 Post-it notes remaining after work. Therefore, the total number of Post-it notes she had after purchasing the package is:

$$80 + 23 = 303$$

Next, we know that Candice used 220 Post-it notes at work. To find the number of Post-it notes in the package, we subtract the number of Post-it notes used from the total number of Post-it notes she had:

$$303 - 220 = 83$$

Thus, the number of Post-it notes in the package is:

83

## RFG (d1 LLaDA)

To determine how many Post-it notes were in the package that Candice purchased, we need to follow these steps:

1. Calculate the total number of Post-it notes Candice used at work. 2. Add the number of Post-it notes she had remaining to the number used at work to find the total number of Post-it notes she had. 3. Subtract the number of Post-it notes she originally had in her purse from the total number of Post-it notes she had to find out how many she purchased in the package.

First, we know that Candice used 220 Post-it notes at work and had 23 Post-it notes remaining. Therefore, the total number of Post-it notes she had is:

$$220 + 23 = 243$$

Next, we know that Candice originally had 80 Post-it notes in her purse. So, the number of Post-it notes she purchased in the package is:

$$243 - 80 = 163$$

Thus, the number of Post-it notes in the package that Candice purchased is:

163

## D.1.4 LLADA 1.5

## Query (GSM8K)

Julia was preparing for a dinner party at her house, where she intended to serve stew. She noticed that she was out of plastic spoons, so she bought a new package of spoons. Later, her husband also bought a package of 5 new spoons and gave them to Julia. While Julia was making the stew, she used three of the spoons to sample her stew. Later, when she went to set the table, she had a total of 12 spoons. How many spoons were in the package that Julia bought?

## Original (LLaDA 1.5)

To determine how many spoons were in the package that Julia bought, we need to follow these steps:

1. **Identify the total number of spoons Julia had before setting the table:** Julia had 12 spoons before setting the table.
2. **Account for the spoons used while making the stew:** Julia used 3 spoons to sample her stew.
3. **Determine the number of spoons given by her husband:** Her husband gave her 5 spoons.
4. **Calculate the total number of spoons before buying the package:** Subtract the spoons used from the total number of spoons before setting the table:

$$12 - 3 = 9$$

So, Julia had 9 spoons before buying the package.

5. **Calculate the number of spoons in the package Julia bought:** Subtract the spoons given by her husband from the total number of spoons before buying the package:

$$9 - 5 = 4$$

Therefore, the package that Julia bought contained 4 spoons.

The final answer is:

4

## RFG (LLaDA 1.5)

To determine how many spoons were in the package that Julia bought, we need to follow these steps:

1. **Identify the total number of spoons Julia had when setting the table:** Julia had 12 spoons when setting the table.
2. **Account for the spoons used while making the stew:** Julia used 3 spoons to sample her stew.
3. **Determine the number of spoons given by her husband:** Her husband gave her 5 spoons.
4. **Calculate the total number of spoons before setting the table:** Let  $x$  be the number of spoons in the package Julia bought. After receiving 5 spoons from her husband, she had  $x + 5$  spoons. After using 3 spoons, she had 12 spoons left.

$$(x + 5) - 3 = 12$$

Simplify the equation:

$$x + 2 = 12$$

Subtract 2 from both sides:

$$x = 10$$

Therefore, Julia bought 10 spoons.

## D.2 MATH-500

## D.2.1 LLADA INSTRUCT

Query (MATH-500)

If  $n \equiv 2 \pmod{7}$ , then find the remainder when  $(n+2)(n+4)(n+6)$  is divided by 7.

Original (LLaDA Instruct)

First, we know that  $n \equiv 2 \pmod{7}$ . We need to find the remainder when  $(n+2)(n+4)(n+6)$  is divided by 7.

Let's compute each term modulo 7:

1.  $n+2 \equiv 2+2 \equiv 4 \pmod{7}$  2.  $n+4 \equiv 2+4 \equiv 6 \pmod{7}$  3.  $n+6 \equiv 2+6 \equiv 0 \pmod{7}$

Now, we multiply these results together:

$$(n+2)(n+4)(n+6) \equiv 4 \cdot 6 \cdot 0 \pmod{7}$$

Since one of the factors is 0, the entire product is:

$$0 \pmod{7}$$

Thus, the remainder when  $(n+2)(n+4)(n+6)$  is divided by 7 is

0

RFG (LLaDA Instruct)

First, we know that  $n \equiv 2 \pmod{7}$ . We need to find the remainder when  $(n+2)(n+4)(n+6)$  is divided by 7.

Let's compute each term modulo 7:

1.  $n+2 \equiv 2+2 \equiv 4 \pmod{7}$  2.  $n+4 \equiv 2+4 \equiv 6 \equiv -1 \pmod{7}$  3.  $n+6 \equiv 2+6 \equiv 8 \equiv 1 \pmod{7}$

Next, we multiply these results together:

$$(4)(-1)(1) \equiv -4 \pmod{7}$$

Finally, we convert  $-4$  to a positive equivalent modulo 7:

$$-4 \equiv 3 \pmod{7}$$

Thus, the remainder when  $(n+2)(n+4)(n+6)$  is divided by 7 is

3

## D.2.2 DREAM INSTRUCT

## Query (MATH-500)

CTwo fair, 6-sided dice are thrown. What is the probability that the product of the two numbers is a multiple of 5? Express your answer as a common fraction.

## Original (Dream Instruct)

To find the probability that the product of the two numbers is a multiple of 5, we need to consider the cases where at least one of the numbers is a multiple of 5. There are 6 possible outcomes for each die, so there are a total of  $6 \times 6 = 36$  possible outcomes when two dice are thrown.

The multiples of 5 on a 6-sided die are 5 itself. Therefore, there is only 1 way for each die to show a multiple of 5. Thus, there are  $1 \times 1 = 1$  outcome where both dice show a multiple of 5.

The probability that the product of the two numbers is a multiple of 5 is the number of favorable outcomes (where at least one die shows a multiple of 5) divided by the total number of outcomes. Since there is only one outcome where both dice show a multiple of 5, the probability is  $\frac{1}{36}$ .

Therefore, the probability that the product of the two numbers is a multiple of 5 is  $\frac{1}{36}$ .

## RFG (Dream Instruct)

To find the probability that the product of the two numbers is a multiple of 5, we need to consider the cases where at least one of the numbers is a multiple of 5. There are 6 possible outcomes for each die, so there are a total of  $6 \times 6 = 36$  possible outcomes when two dice are thrown.

The multiples of 5 on a 6-sided die are 5 itself. Therefore, there is only one way for each die to show a multiple of 5. The probability of neither die showing a multiple of 5 is  $(\frac{5}{6})^2 = \frac{25}{36}$ . Thus, the probability that at least one die shows a multiple of 5 is  $1 - \frac{25}{36} = \frac{11}{36}$ .

Therefore, the probability that the product of the two numbers is a multiple of 5 is  $\frac{11}{36}$ .

## D.2.3 D1-LLADA

## Query (MATH-500)

Compute  $\arcsin(-\frac{1}{2})$ . Express your answer in radians.

## Original (d1 LLaDA)

To compute  $\arcsin(-\frac{1}{2})$ , we need to find the angle  $\theta$  such that  $\sin(\theta) = -\frac{1}{2}$ . The angle whose sine is  $-\frac{1}{2}$  is  $\frac{5\pi}{6}$  radians.

Thus, the answer is:

$$\frac{5\pi}{6}$$

## RFG (d1 LLaDA)

To compute  $\arcsin(-\frac{1}{2})$ , we need to find the angle  $\theta$  such that  $\sin \theta = -\frac{1}{2}$  and  $\theta$  lies within the range of the arcsine function, which is  $[-\frac{\pi}{2}, \frac{\pi}{2}]$ .

We know that  $\sin(-\frac{\pi}{6}) = -\frac{1}{2}$ . The angle  $-\frac{\pi}{6}$  is within the range  $[-\frac{\pi}{2}, \frac{\pi}{2}]$ .

Thus,  $\arcsin(-\frac{1}{2}) = -\frac{\pi}{6}$ .

$$-\frac{\pi}{6}$$

## D.2.4 LLADA 1.5

Query (MATH-500)

What power of 4 is equal to 8? Express your answer as a common fraction.

Original (LLaDA 1.5)

To determine what power of 4 is equal to 8, we need to solve the equation  $4^x = 8$ .  
First, we can express 8 as a power of 4:

$$8 = 4^2$$

So, the equation becomes:

$$4^x = 4^2$$

Since the bases are the same, we can set the exponents equal to each other:

$$x = 2$$

Therefore, the power of 4 that is equal to 8 is:

$$\boxed{2}$$

RFG (LLaDA 1.5)

To determine what power of 4 is equal to 8, we need to solve the equation  $4^x = 8$ .  
First, let's express 8 as a power of 2:

$$8 = 2^3$$

So, we can rewrite the equation as:

$$4^x = 2^3$$

Next, we express 4 as a power of 2:

$$4 = 2^2$$

Thus, we can rewrite the equation as:

$$(2^2)^x = 2^3$$

Using the power of a power property  $(a^m)^n = a^{mn}$ , we get:

$$2^{2x} = 2^3$$

Since the bases are the same, we can set the exponents equal to each other:

$$2x = 3$$

Solving for  $x$ , we divide both sides by 2:

$$x = \frac{3}{2}$$

Therefore, the power of 4 that is equal to 8 is:

$$\boxed{\frac{3}{2}}$$

### D.3 HUMANEVAL

#### D.3.1 LLADA INSTRUCT

##### Query (HumanEval)

```
def sort_array(array):
    """
    Given an array of non-negative integers, return a copy of the given array after sorting,
    you will sort the given array in ascending order if the sum
    ( first index value, last index value) is odd,
    or sort it in descending order if the sum( first index value, last index value) is even.

    Note:
    * don't change the given array.

    Examples:
    * sort_array([]) => []
    * sort_array([5]) => [5]
    * sort_array([2, 4, 3, 0, 1, 5]) => [0, 1, 2, 3, 4, 5]
    * sort_array([2, 4, 3, 0, 1, 5, 6]) => [6, 5, 4, 3, 2, 1, 0]
    """
```

##### Original (LLaDA Instruct)

```
if len(array) == 0:
    return []

sum_last = array[0] + array[-1]
if sum_last % 2 == 1:
    return sorted(array, reverse=True)
else:
    return sorted(array, reverse=False)
```

##### RFG (LLaDA Instruct)

```
if len(array) == 0:
    return []

sum_last = array[0] + array[-1]

if sum_last % 2 == 1:
    return sorted(array)
else:
    return sorted(array, reverse=True)
```

#### D.3.2 DREAM INSTRUCT

##### Query (HumanEval)

```
def sort_numbers(numbers: str) -> str:
    """ Input is a space-delimited string of numerals from 'zero' to 'nine'.
    Valid choices are 'zero', 'one', 'two', 'three', 'four', 'five', 'six',
    'seven', 'eight' and 'nine'.
    Return the string with numbers sorted from smallest to largest
    >>> sort_numbers('three one five')
    'one three five'
    """
```

##### Original (Dream Instruct)

```
numbers = numbers.split()
numbers.sort(key=lambda x: numbers.index(x))
return ' '.join(numbers)
```

##### RFG (Dream Instruct)

```
valid_numbers = ['zero', 'one', 'two', 'three', 'four', 'five', 'six', 'seven', 'eight', 'nine']
numbers = [valid_numbers.index(num) for num in numbers.split()]
numbers.sort()
return ' '.join(str(valid_numbers[num]) for num in numbers)
```

### D.3.3 DIFFUCODER

#### Query (HumanEval)

```
def how_many_times(string: str, substring: str) -> int:
    """ Find how many times a given substring can be found in the original string.
    Count overlaping cases.
    >>> how_many_times('', 'a')
    0
    >>> how_many_times('aaa', 'a')
    3
    >>> how_many_times('aaaa', 'aa')
    3
    """
    """
```

#### Original (DiffuCoder)

```
return string.count(substring)
```

#### RFG (DiffuCoder)

```
count = 0
for i in range(len(string) - len(substring) + 1):
    if string[i:i+len(substring)] == substring:
        count += 1
return count
```

## D.4 MBPP

### D.4.1 LLaDA INSTRUCT

#### Query (MBPP)

Write a function to append the given list to the given tuples. Your code should pass the test:

```
assert add_lists([5, 6, 7], (9, 10)) == (9, 10, 5, 6, 7)
```

#### Original (LLaDA Instruct)

```
def add_lists(lst, tuples):
    return tuple(lst + tuples)
```

#### RFG (LLaDA Instruct)

```
def add_lists(lst, tuples):
    # Convert the tuples to a list
    tuples_list = list(tuples)
    # Append the given list to the tuples list
    tuples_list.extend(lst)
    # Convert the list back to a tuple
    return tuple(tuples_list)
```

## D.4.2 DREAM INSTRUCT

## Query (MBPP)

Write a python function to find the number of numbers with an odd value when rotating a binary string the given number of times. Your code should pass the test:

```
assert odd_Equivalent("011001",6) == 3
```

## Original (Dream Instruct)

```
def odd_Equivalent(binary_string, rotations):
    n = len(binary_string)
    odd_count = 0
    for i in range(rotations):
        binary_string = binary_string[1:] + binary_string[0]
        odd_count += sum(int(bit) for bit in binary_string) % 2
    return odd_count
```

## RFG (Dream Instruct)

```
def odd_Equivalent(binary, n):
    count = 0
    for i in range(n):
        binary = binary[1:] + binary[0]
        if int(binary, 2) % 2 == 1:
            count += 1
    return count
```

## D.4.3 DIFFUCODER

## Query (MBPP)

Write a function to find the largest sum of a contiguous array in the modified array which is formed by repeating the given array k times. Your code should pass the test:

```
assert max_sub_array_sum_repeated([10, 20, -30, -1], 4, 3) == 30
```

## Original (DiffuCoder)

```
def max_sub_array_sum_repeated(arr, n, k):
    repeated_arr = arr * k
    max_sum = float('-inf')
    for i in range(n):
        current_sum = 0
        for j in range(i, i + n):
            current_sum += repeated_arr[j % n]
            max_sum = max(max_sum, current_sum)
    return max_sum
```

## RFG (DiffuCoder)

```
def max_sub_array_sum_repeated(arr, n, k):
    repeated_arr = arr * k
    max_sum = float('-inf')
    for i in range(n * k):
        current_sum = 0
        for j in range(i, n * k):
            current_sum += repeated_arr[j % n]
            max_sum = max(max_sum, current_sum)
    return max_sum
```