

# Reward-Weighted Sampling: Enhancing Non-Autoregressive Characteristics in Masked Diffusion LLMs

Daehoon Gwak<sup>1\*</sup> Minseo Jung<sup>2\*</sup> Junwoo Park<sup>1</sup> Minho Park<sup>1</sup>  
 ChaeHun Park<sup>1</sup> Junha Hyung<sup>1</sup> Jaegul Choo<sup>1</sup>

<sup>1</sup>KAIST AI

<sup>2</sup>Applied Artificial Intelligence, Sungkyunkwan University  
 {daehoon.gwak, jchoo}@kaist.ac.kr jms020123@g.skku.edu

## Abstract

Masked diffusion models (MDMs) offer a promising non-autoregressive alternative for large language modeling. Standard decoding methods for MDMs, such as confidence-based sampling, select tokens independently based on individual token confidences at each diffusion step. However, we observe that this independent token selection often results in generation orders resembling sequential autoregressive processes, limiting the advantages of non-autoregressive modeling. To mitigate this pheonomenon, we propose Reward-Weighted Sampling (RWS), a novel decoding strategy that leverages an external reward model to provide a principled global signal during the iterative diffusion process. Specifically, at each diffusion step, RWS evaluates the quality of the entire intermediate sequence and scales token logits accordingly, guiding token selection by integrating global sequence-level coherence. This method selectively increases the confidence of tokens that initially have lower scores, thereby promoting a more non-autoregressive generation order. Furthermore, we provide theoretical justification showing that reward-weighted logit scaling induces beneficial rank reversals in token selection and consistently improves expected reward. Experiments demonstrate that RWS significantly promotes non-autoregressive generation orders, leading to improvements across multiple evaluation metrics. These results highlight the effectiveness of integrating global signals in enhancing both the non-autoregressive properties and overall performance of MDMs.

## 1 Introduction

Large Language Models (LLMs) have achieved remarkable success across diverse natural language tasks, predominantly through autoregressive models (ARMs) (Brown et al., 2020; Lewis et al., 2020), which generate text sequentially, token by token.

\*Equal contribution.

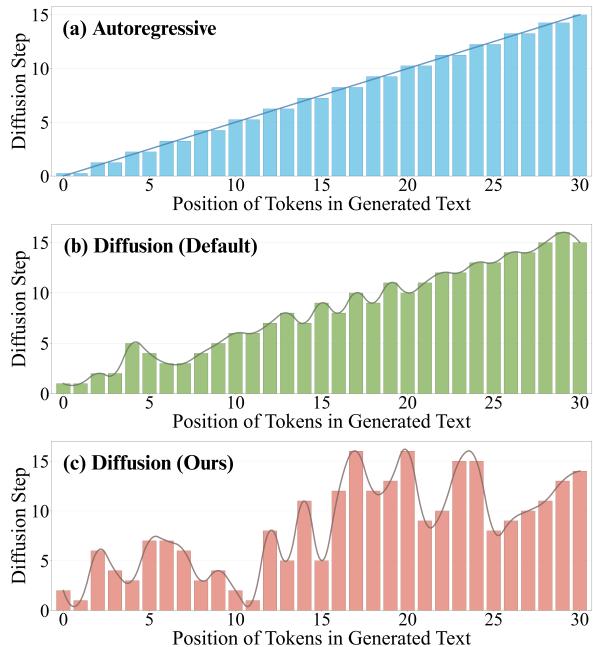


Figure 1: Representative example illustrating token selection order for (a) Autoregressive (top), (b) Diffusion LLM with confidence-based sampling (middle), and (c) our proposed Reward-Weighted Sampling (RWS) method (bottom). The x-axis shows token positions within the generated sequence, and the y-axis shows the diffusion step at which two tokens are selected. Confidence-based sampling closely follows the sequential autoregressive pattern, while RWS promotes a distinctly non-autoregressive selection pattern.

Recently, large-scale Masked Diffusion Models (MDMs), such as LLaDA (Nie et al., 2025), have emerged as promising non-autoregressive alternatives. Unlike ARMs, MDMs iteratively unmask tokens in parallel, leveraging full bidirectional context. This approach can mitigate cumulative errors commonly seen in autoregressive generation, where early prediction errors propagate through subsequent tokens, often causing semantic drift and reduced coherence in longer sequences (Bengio et al., 2015; Ranzato et al., 2016).

Despite their inherent non-autoregressive potential, decoding methods specifically tailored to

MDMs are relatively under-explored. Standard decoding approaches for MDMs, such as confidence-based sampling (Nie et al., 2025), rely solely on individual token confidences at each diffusion step, selecting tokens with the highest confidence score to unmask. As illustrated in Figure 1, this token selection method often produces generation orders closely resembling sequential autoregressive decoding, represented by a diagonal selection pattern. This phenomenon could be interpreted as a position bias effect commonly observed in masked language models (Wang and Cho, 2019; Ghazvininejad et al., 2019), where masked tokens adjacent to already unmasked tokens tend to receive higher confidence scores. Consequently, the model often unmasks tokens in a locally sequential manner, potentially limiting its ability to fully exploit global, bidirectional context and impacting tasks that require global coherence or logical consistency.

To mitigate this limitation, we propose Reward-Weighted Sampling (RWS), a novel decoding method explicitly designed to exploit the non-autoregressive characteristics of MDMs by introducing a global preference signal. At each diffusion step, RWS evaluates the quality of intermediate sequences generated by the model and computes a global reward signal using an external reward model. This global reward reflects the overall coherence and quality of the intermediate sequence, beyond the local token-level confidences. The logits of individual tokens are then scaled according to this reward, selectively boosting tokens initially assigned lower local confidence scores, thus promoting a more diversified, non-autoregressive generation order as illustrated in Figure 1.

Furthermore, we provide rigorous theoretical analysis demonstrating how reward-weighted logit scaling influences token selection. Specifically, we derive precise conditions under which multiplicative logit scaling induces beneficial rank reversals in token probabilities (Theorem 1), thereby promoting more flexible generation patterns. Additionally, we prove that reward-based logit scaling guarantees monotonically improved expected reward at each decoding step (Theorem 2), providing theoretical justification for the effectiveness of our method in exploiting non-autoregressive decoding to enhance sequence generation quality.

In summary, this paper makes the following contributions:

- We introduce Reward-Weighted Sampling, a

decoding method that integrates global reward signals to systematically encourage non-autoregressive token selection.

- We provide theoretical analyses proving that logit scaling causes principled rank reversals in token selection and guarantees improvement of expected rewards during decoding, establishing the mathematical foundation for enhanced non-autoregressive behavior.
- Extensive experiments across diverse tasks demonstrate that our method consistently improves generation quality compared to baseline methods, while showing a strong correlation between non-sequential generation, measured by our proposed metric, Generation Order Deviation (GOD).

Our findings indicate that integrating a meaningful global reward signal effectively enhances the non-autoregressive capabilities and overall generation quality of diffusion-based language models.

## 2 Background: Standard Sampling for Masked Diffusion Models

In this section, we introduce the standard decoding methods used in MDMs, highlighting the limitations we aim to address. For a detailed explanation of the underlying MDM framework itself, we refer readers to Appendix A.

Generating text from a pre-trained MDM, such as generating a response  $r_0$  conditioned on a prompt  $p_0$ , involves simulating the reverse diffusion process in discrete steps. Let the target generation length be  $L'$ . The process starts with an initial sequence  $x^{(T)} = [p_0, [\text{MASK}], \dots, [\text{MASK}]]$ , where the prompt  $p_0$  is followed by  $L'$  mask tokens, denoted as  $[\text{MASK}]$ . The goal is to iteratively denoise this sequence over  $T$  discrete steps (indexed  $t = T, T-1, \dots, 1$ ) to obtain the final sequence  $x^{(0)} \approx [p_0, r_0]$ .

The standard iterative sampling procedure at step  $t$  (transitioning from state  $x^{(t)}$  to  $x^{(t-1)}$ ) typically involves three steps. **Step 1 (Prediction):** The current partially masked sequence  $x^{(t)}$  is fed into the trained mask predictor  $p_\theta$  to obtain logits for all masked positions. Let  $\hat{x}_0$  be the sequence obtained by replacing each masked position with the most likely original token predicted from these logits (e.g., via argmax). **Step 2 (Token Selection):** A subset of masked positions in  $x^{(t)}$  is selected for updating, typically the top  $k_t$  positions with

the highest confidence, ensuring gradual denoising (e.g., approximately  $L'/T$  tokens per step). **Step 3 (Update):** The next state  $x^{(t-1)}$  is formed by replacing the selected  $k_t$  masked positions in  $x^{(t)}$  with their corresponding predicted tokens from  $\hat{x}_0$ .

A common strategy for Step 2, employed by LLaDA and forming the baseline for our work, is **Confidence-Based Sampling**. In this approach, token confidence is computed from softmax probabilities  $p_\theta(\hat{x}_0 | x^{(t)})$ , the top  $k_t$  tokens with highest confidence are selected, and these tokens are unmasked in  $x^{(t-1)}$  while the rest remain masked for subsequent refinement.

It is important to note that prior studies (Ghazvininejad et al., 2019; Wang and Cho, 2019) have demonstrated that in masked language models, prediction confidence is heavily influenced by surrounding context. Specifically, masked tokens directly adjacent to already unmasked tokens—such as prompt tokens—tend to receive higher confidence scores due to their immediate contextual proximity. As a result, tokens closer to previously unmasked content are generally selected first, causing a positional bias in which the model tends to unmask tokens sequentially from left to right, despite its inherent bidirectional architecture. This positional bias unintentionally promotes autoregressive-like decoding behavior, potentially limiting the model’s ability to leverage the full bidirectional context effectively.<sup>1</sup>

### 3 Methodology: Reward-Weighted Sampling (RWS)

#### 3.1 Overview

Reward-Weighted Sampling (RWS) explicitly guides the token selection process during MDM decoding by incorporating feedback from an external reward model. To achieve this, RWS introduces a global reward signal into the iterative decoding process. Multiplicative scaling of logits provides an intuitive mechanism for incorporating such global signals, as scaling logits directly adjusts the differences between token probabilities due to the exponential nature of the softmax function (Hinton et al., 2015; Holtzman et al., 2020). Specifically, increasing logits amplifies probability gaps, thereby reinforcing the selection of high-confidence tokens, whereas decreasing logits reduces these gaps, enabling lower-confidence tokens to become compar-

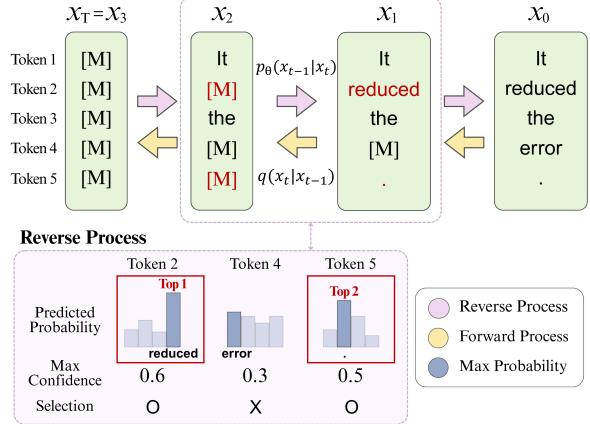


Figure 2: Illustration of the Masked Diffusion Model (MDM) forward (masking) and reverse (unmasking) processes. Starting from the fully masked sequence ( $x_T$ ), the reverse process progressively reveals tokens using predicted probabilities at each diffusion step. The highlighted section illustrates token selection between steps  $x_2$  and  $x_1$ : tokens with the highest prediction confidence are unmasked first (tokens 2 and 5), while lower-confidence tokens (token 4) remain masked.

atively more likely. Utilizing an external reward model to guide this scaling ensures that token selection adjustments meaningfully reflect the quality of candidate sequences, thus avoiding arbitrary or potentially harmful modifications.

#### 3.2 Reward-Weighted Sampling Algorithm

RWS modifies the standard iterative MDM sampling process described in Section 2. Starting from the initial state  $x^{(T)} = [p_0, [\text{MASK}], \dots, [\text{MASK}]]$  (prompt  $p_0$  followed by  $L'$  mask tokens), RWS proceeds iteratively for  $T$  steps ( $t = T, T-1, \dots, 1$ ) to produce the final sequence  $x^{(0)}$ . Specifically, at each diffusion step  $t$  (going from masked sequence  $x^{(t)}$  to the next state  $x^{(t-1)}$ ), RWS follows these four steps:

**Step 1: Potential Full Sequence Prediction.** At the current step  $t$  with state  $x^{(t)}$ , we first generate a full candidate sequence sequence  $\tilde{x}_0^{(t)}$ . Specifically, we obtain logits  $L^{(t)}$  from the mask predictor  $p_\theta$  given the current masked sequence  $x^{(t)}$ , and then greedily predict the most likely token (argmax) for all currently masked positions. Let  $M^{(t)}$  be the set of indices corresponding to masked positions in  $x^{(t)}$ . Then,  $\tilde{x}_0^{(t)}$  is formed by setting  $\tilde{x}_0^j = \text{argmax}_{\text{token}} L_{j,\text{token}}^{(t)}$  for  $j \in M^{(t)}$ , and  $\tilde{x}_0^j = x^{(t),j}$  for  $j \notin M^{(t)}$ . This  $\tilde{x}_0^{(t)}$  represents a possible completed output given the current masked context  $x^{(t)}$ .

**Step 2: Reward Evaluation and Normalization.** The generated part (response  $r_0$ ) within the

<sup>1</sup>For a broader overview of diffusion language models and other related works, please refer to Appendix B.

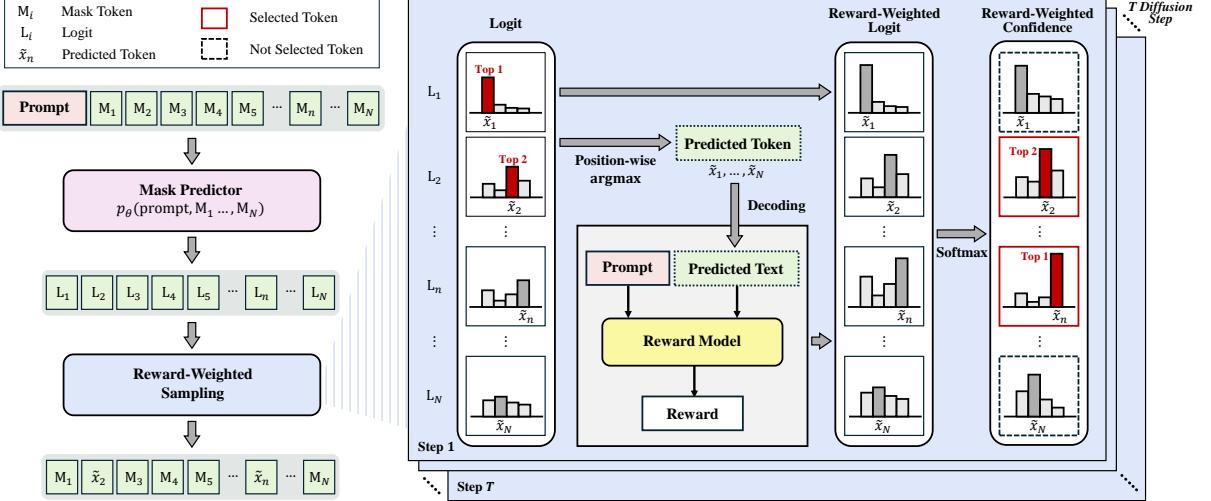


Figure 3: Overview of Reward-Weighted Sampling (RWS) method. At each diffusion step, the model predicts candidate tokens based on local confidence scores. An external reward model then evaluates the global quality of the predicted sequence, providing a reward-based signal. This signal adaptively scales logits, adjusting token selection priorities to promote non-autoregressive generation patterns and improve overall sequence coherence.

potential full sequence  $\tilde{x}_0^{(t)}$  is extracted and decoded into text. This text, along with the original prompt  $p_0$ , is evaluated by a external reward model  $R(p_0, r_0)$  to obtain a raw reward score  $r_{\text{raw}}^{(t)}$ . As shown in Appendix C, due to significant variations in reward statistics across different reward models, we normalize the raw reward using pre-computed mean  $\mu_R$  and standard deviation  $\sigma_R$  from the Nec-tar validation set (Zhu et al., 2023), computed as  $r^{(t)} = (r_{\text{raw}}^{(t)} - \mu_R)/\sigma_R$ . This normalized reward  $r^{(t)}$  reflects the relative quality of the potential completion compared to typical outputs.

**Step 3: Reward-Weighted Logit Scaling.** This is the core step where the reward signal directly influences token selection priorities by modifying the original logits  $L^{(t)}$  obtained from  $p_\theta(x^{(t)})$  using the normalized reward  $r^{(t)}$ . Specifically, the modified logits  $L'^{(t)}$  are calculated as:

$$L'_{j,\text{token}}^{(t)} = L_{j,\text{token}}^{(t)} \times s_R \times \sqrt{\sigma(r^{(t)})} + \epsilon$$

for all positions  $j$  and tokens. Here,  $s_R$  is the reward scale hyperparameter controlling the guidance strength,  $\sigma(\cdot)$  is the sigmoid function, and  $\epsilon$  is a small constant (e.g.,  $10^{-5}$ ) for numerical stability. This scaling adjusts the original logits according to the reward signal, influencing the likelihood of token selection. We employ the sigmoid function to normalize the reward into a stable range between 0 and 1. The square root is subsequently applied to prevent excessively small scaling values for negative rewards.

Importantly, as shown in our theoretical analy-

sis (Section 4), larger multiplicative scaling values increase the likelihood of rank reversals in token selection, promoting more non-autoregressive generation patterns. By linking this scaling factor to the current sequence’s reward, our method adaptively encourages more rank reversals when the intermediate sequence shows promising quality (higher reward), and fewer when the quality is poor (lower reward). This adaptive approach ensures that the degree of non-autoregressive behavior is dynamically adjusted based on generation quality.

**Step 4: Guided Token Selection.** Finally, the selection of  $k_t$  tokens to unmask proceeds similarly to the standard Confidence-Based Sampling (Section 2), but critically, it uses the *modified* logits  $L'^{(t)}$  to compute confidence scores. Typically, the confidence is derived from the softmax probabilities of the argmax predictions based on  $L'^{(t)}$ . The  $k_t$  masked positions with the highest confidence derived from the reward-adjusted logits  $L'^{(t)}$  are selected. The corresponding predicted tokens replace the mask tokens at these selected positions to form the next state  $x^{(t-1)}$ . Note that while the rank order of confidence scores between each token position can change due to reward scaling, the argmax token predictions at each position remain the same regardless of reward scaling.

Repeating these steps for all diffusion steps, RWS leverages external rewards to produce more effective, non-autoregressive text generation.

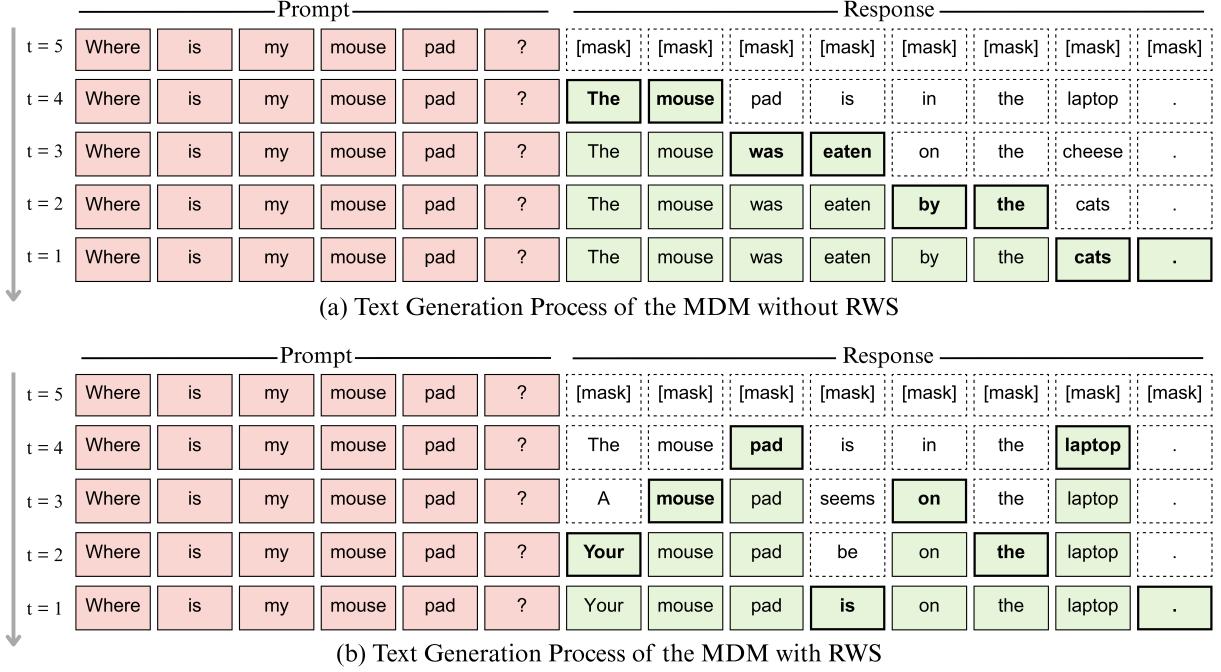


Figure 4: Illustration of how RWS promotes non-autoregressive token selection. Decoding proceeds from  $t = 5$  (fully masked) down to  $t = 1$  (fully unmasked). (a) Default sampling unmasks tokens sequentially, prematurely fixing earlier tokens ("The mouse") and resulting in an incorrect, contextually mismatched output ("The mouse was eaten by the cats."). (b) RWS enables flexible, non-sequential token selection, correctly identifying tokens based on global context, generating a coherent output ("Your mouse pad is on the laptop.").

## 4 Theoretical Analysis

This section provides a theoretical analysis illustrating how scaling logits multiplicatively can change which token is most likely.

### 4.1 Setup and Notation

Let  $a, b \in \mathbb{R}^d$  be two logit vectors obtained at some decoding step.<sup>2</sup> For any vector  $x \in \mathbb{R}^d$  and scaling factor  $r > 0$ , define the scaled softmax distribution:

$$q_i^x(r) = \frac{\exp(rx_i)}{\sum_{j=1}^d \exp(rx_j)}.$$

Let  $p_i^x = q_i^x(1)$  denote the original softmax probabilities, and write  $P^x(r) = \max_i q_i^x(r)$ . Throughout, we assume:

$$\max_i p_i^a > \max_i p_i^b. \quad (1)$$

### 4.2 Intuition: Why Scaling Can Change Token Selection

Scaling logits by a factor  $r$  uniformly shrinks or expands all logit gaps  $\Delta_k = l_{\max} - l_k$ :

- $r < 1$ : All gaps shrink, flattening the probability distribution. Here, the total sum of gaps

<sup>2</sup>In practice,  $a$  and  $b$  could come from two different masked positions within the same MDM step (§2).

determines which token's maximum probability decreases faster.

- $r > 1$ : All gaps expand, sharpening the distribution. In this case, the smallest gap (the closest competing token) determines which token's maximum probability approaches 1 more slowly.

Thus, a reversal in token selection occurs if (a) the currently highest-ranked token vector has a larger total gap sum or (b) a smaller nearest-competitor gap compared to the alternative. The theorem below formalizes this concept.

### 4.3 Theoretical Conditions for Rank Reversal

Write the gaps of  $a$  and  $b$  relative to their top indices  $t_a, t_b$  as

$$\Delta_k^a = a_{t_a} - a_k, \quad \Delta_k^b = b_{t_b} - b_k \quad (k \neq t_a, t_b).$$

Define the total and minimal gaps

$$\begin{aligned} \Sigma_a &= \sum_{k \neq t_a} \Delta_k^a, & \delta_a &= \min_{k \neq t_a} \Delta_k^a, \\ \Sigma_b &= \sum_{k \neq t_b} \Delta_k^b, & \delta_b &= \min_{k \neq t_b} \Delta_k^b. \end{aligned} \quad (2)$$

**Theorem 1** (Rank Reversal Conditions). *Under Eq. (1), exactly one of the following holds:*

Reward Model	GOD ( $\uparrow$ )	RWS Win (%)	Draw (%)	Default Win (%)
None (Default)	1.426	-	-	-
Skywork	<b>2.607</b>	<b>60.2</b>	3.8	36.0
LlamaOB	<b>2.621</b>	<b>60.9</b>	4.4	34.7
Intern	<b>2.599</b>	<b>58.7</b>	3.6	37.8
Eurus	<b>2.658</b>	<b>58.7</b>	3.6	37.8
GRM	<b>2.621</b>	<b>55.8</b>	4.0	40.2
QRM	<b>2.394</b>	<b>60.2</b>	3.3	36.4

Table 1: Generation Order Deviation (GOD) and win-rate comparison between Reward-Weighted Sampling (RWS) and the confidence-based baseline (Default) across different automated reward models. Higher GOD values indicate improved non-autoregressive behavior.

- (a) If  $\Sigma_a < \Sigma_b$ , there exists a unique  $r^* \in (0, 1)$  such that  $P^a(r^*) = P^b(r^*)$  and  $P^a(r) < P^b(r)$  for all  $0 < r < r^*$ .
- (b) If  $\delta_a > \delta_b$ , there exists a unique  $r^\dagger > 1$  such that  $P^a(r^\dagger) = P^b(r^\dagger)$  and  $P^a(r) < P^b(r)$  for all  $r > r^\dagger$ .
- (c) Otherwise  $P^a(r) > P^b(r)$  for every  $r > 0$ ; no reversal occurs.

*Proof.* See Appendix D for the detailed proof.  $\square$

**Theorem 2** (Reward-Monotonic Guidance). *Let  $p_\theta(x) \propto \exp L_\theta(x)$  be the base decoder and  $p_r(x) \propto p_\theta(x) e^{rR(x)}$  its reward-scaled variant with any  $r > 0$ . Then*

$$\mathbb{E}_{p_r}[R] \geq \mathbb{E}_{p_\theta}[R], \quad (3)$$

with strict inequality if  $\text{Var}_{p_r}[R] > 0$ .

*Proof.* The full proof is provided in Appendix E.  $\square$

**Corollary 1** (Per-step Improvement in RWS). *At diffusion step  $t$  of the reverse process, let  $p_\theta(\cdot | x^{(t)})$  be the model distribution conditioned on the current mask and let  $r^{(t)} > 0$  be the scale produced from the (normalised) reward. Sampling the token update  $x_{\text{RWS}}^{(t)}$  from  $p_{r^{(t)}}$  guarantees*

$$\mathbb{E}[R(x_{\text{RWS}}^{(t)})] \geq \mathbb{E}[R(x^{(t)})],$$

*with equality only if  $R$  is constant on the support of  $p_{r^{(t)}}$ .*

*Proof.* Given the current masked sequence  $x^{(t)}$ , the distributions  $p_\theta(\cdot | x^{(t)})$  and  $p_{r^{(t)}}(\cdot | x^{(t)})$  correspond exactly to the setting of Theorem 2. Thus, directly applying Theorem 2, we obtain:

$$\mathbb{E}_{p_{r^{(t)}}(\cdot | x^{(t)})}[R(x_{\text{RWS}}^{(t)})] \geq \mathbb{E}_{p_\theta(\cdot | x^{(t)})}[R(x^{(t)})],$$

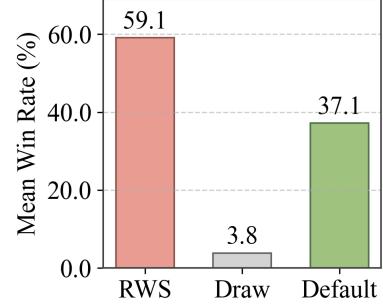


Figure 5: Average win rate comparison across all six reward models. RWS outperforms the Default method by a significant margin.

with strict inequality if  $R$  is not constant on the support of  $p_{r^{(t)}}(\cdot | x^{(t)})$ . Since this inequality holds for any arbitrary masked sequence  $x^{(t)}$ , it directly implies the stated per-step improvement of the Reward-Weighted Sampling approach.  $\square$

**Interpretation.** Theorems 1 and 2 explain why Reward-Weighted Sampling (RWS) improves the decoding process. Specifically, they show that scaling logits based on a reward model consistently increases the expected reward of the generated tokens. Additionally, these theorems identify clear conditions under which tokens with initially lower confidence become preferred, effectively promoting a non-autoregressive generation pattern.

Figure 4 provides a concrete illustration of this theoretical insight: under standard sampling, the diffusion model follows a locally autoregressive pattern, resulting in a contextually incorrect generation. In contrast, RWS induces beneficial token selection rank reversals, producing a coherent and globally consistent output. In short, RWS systematically selects better tokens by combining local confidence with global reward signals.

**Comparison with Temperature Scaling.** We also provide a detailed theoretical and empirical comparison between RWS and the widely used softmax temperature scaling method in Appendix G.

## 5 Experiments

### 5.1 Experimental Setup

We use LLaDA-8B-Instruct as our base diffusion LLM across all experiments. We evaluate our method using six publicly available reward models from RewardBench (Lambert et al., 2024), summarized in Appendix H. We conduct experiments on three distinct benchmarks: RewardBench (Lambert et al., 2024), MT-Bench (Zheng et al., 2023), and a

keyword-constrained generation task designed to test non-autoregressive capabilities.

We evaluate model performance using our proposed metric, *Generation Order Deviation (GOD)*, which quantifies how far the actual token generation order deviates from a strictly left-to-right (autoregressive) decoding sequence. Formally, GOD is defined as:

$$\text{GOD} = \frac{1}{n} \sum_{i=1}^n |a_i - o_i|,$$

where  $n$  is the total number of tokens in the generated sequence,  $o_i$  is the position at which the  $i$ -th token was actually decoded by the diffusion model, and  $a_i$  is the position at which the  $i$ -th token would have been decoded under strictly autoregressive (left-to-right) decoding. Higher GOD indicates stronger non-autoregressive behavior; for instance, a GOD of 2.1 implies tokens are decoded, on average, 2.1 positions away from the standard autoregressive order.

Additional standard metrics such as *Win Rate* and *Perplexity* are also employed. Further details on reward models, datasets, and general implementation details can be found in Appendix I.

## 5.2 Reward-Weighted vs. Standard Sampling

We first evaluate the effectiveness of RWS compared to the default confidence-based sampling method across different reward models. As shown in Table 1 and Figure 5, RWS consistently outperforms the default method across all six reward models, with win rates ranging from 55.8% to 60.9%.

Notably, the consistent improvement across multiple reward models with different architectures and training objectives demonstrates the robustness of our approach. The average win rate of approximately 59% represents a non-trivial improvement over baseline methods (37.1%), particularly considering that these evaluations were conducted on the RewardBench dataset, which was specifically designed to assess how well reward models align with human preferences.

The results indicate that incorporating global reward signals during the diffusion process effectively guides the model toward higher-quality generations. This improvement can be attributed to RWS’s ability to deviate from the sequential left-to-right generation pattern inherent in standard confidence-based sampling, instead promoting a more flexible generation order that better utilizes

Reward Model	RWS Win (%)	Default Win (%)
Skywork	<b>65.0</b>	35.0
LlamaOB	<b>62.5</b>	37.5
Intern	<b>60.6</b>	39.4
Eurus	<b>61.9</b>	38.1
GRM	<b>63.1</b>	36.9
QRM	<b>60.0</b>	40.0

Table 2: LLM judge-based win rates comparing RWS against the default method, averaged over two runs with reversed answer ordering to control for position bias.

the non-autoregressive capabilities of MDMs. This is quantitatively supported by the GOD metric as demonstrated in Table 1, where RWS achieves significantly higher values compared to the default method, indicating a greater deviation from strict left-to-right generation.

## 5.3 Performance on Multi-Turn Response Tasks

To evaluate the effectiveness of RWS in more complex conversational settings, we conduct experiments on multi-turn dialogue prompts from MT-Bench. For this experiment, we employ LLM-as-a-judge evaluation methodology, following Zheng et al. (2023) and Lin and Chen (2023). Further details on the exact evaluation prompt and setup are provided in Appendix J.

Table 2 presents the results of this evaluation, comparing RWS against the default confidence-based sampling method. To control for potential position bias in the judging process (Zheng et al., 2023), we conduct evaluations in two configurations, one with RWS outputs presented as choice 1 and another with RWS outputs as choice 2, and reported the averaged win rates.

The results show that RWS consistently outperforms the default method across all reward models, with win rates ranging from 60.0% to 65.0%. This consistent advantage regardless of response order confirms that the improvements are genuine and not artifacts of evaluation bias. The full results for each ordering configuration are provided in Appendix K.

We further provide a qualitative example illustrating how RWS improves coherence in multi-turn settings. A representative example comparing RWS and Default Sampling across a two-turn dialogue is provided in Appendix L. While Default Sampling generates repetitive and vague responses, RWS consistently produces coherent and contextu-

Prompt	Anchor Keywords
Create a brief analysis of climate change using all anchor keywords. Make sure the text reads naturally and maintains logical flow.	<i>global temperatures, carbon emissions, international agreements, renewable technologies</i>
Write a short business strategy description incorporating all anchor keywords. The text should be coherent and professional.	<i>market analysis, competitive advantage, revenue streams, strategic partnerships</i>

Table 3: Examples from keyword-constrained generation dataset. Each input provided to the model consists of a prompt along with specific anchor keywords that must be incorporated within the generated response.

ally rich outputs, clearly demonstrating enhanced cross-turn reasoning and coherence.

The multi-turn setting is particularly challenging as it requires maintaining coherence across longer contexts and adhering to conversational flow—aspects that benefit from non-autoregressive generation’s ability to consider global context. RWS’s strong performance in this setting demonstrates its effectiveness in leveraging bidirectional information for complex reasoning and contextual understanding.

#### 5.4 Performance on Keyword-Constraint Generation Tasks

Our third experiment evaluates model performance on a keyword-constrained generation task, specifically designed to test non-autoregressive capabilities. Inspired by the CommonGen dataset (Lin et al., 2020), which requires models to generate coherent sentences incorporating a given set of concepts, we constructed our dataset to further increase task difficulty by requiring the incorporation of predetermined anchor keywords into longer, contextually coherent responses. This setup explicitly challenges the model’s ability to maintain global coherence and natural fluency. Details of dataset creation, including keyword selection criteria and validation processes, are provided in Appendix M. Table 3 shows representative examples from our keyword-constrained dataset.

For evaluation, we measure keyword inclusion rate, text perplexity, and GOD. While keyword inclusion directly evaluates the model’s ability to satisfy task constraints, perplexity is used as a widely-accepted metric to quantitatively assess fluency and coherence in generated text. GOD further quantifies the non-autoregressive characteristics of the generation.

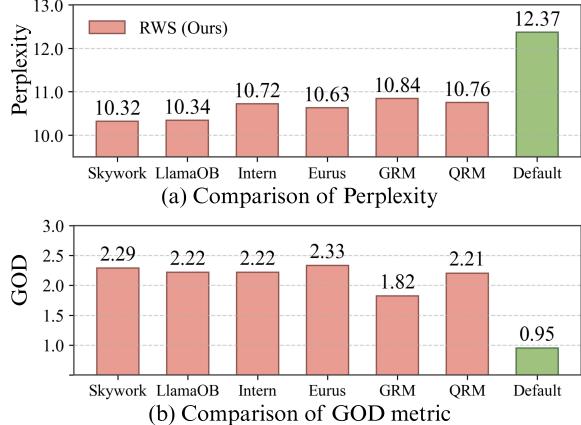


Figure 6: Perplexity (top) and GOD (bottom) results on keyword-constrained generation tasks. Lower perplexity indicates better fluency, while higher GOD reflects stronger non-autoregressive behavior.

While all methods achieved similar keyword inclusion rates (3.6–3.7 keywords out of 4), significant differences emerged in perplexity scores and GOD, as shown in Figure 6. RWS consistently achieves lower perplexity across all reward model configurations, with improvements ranging from 12% to 16% compared to the default method. Lower perplexity indicates more fluent and natural text, suggesting that RWS produces more coherent outputs while still satisfying the keyword constraints.

Additionally, RWS achieves significantly higher GOD values (on average more than double the GOD values of the default method), clearly demonstrating that our method successfully promotes non-autoregressive token selection patterns, effectively leveraging the diffusion model’s bidirectional context to improve overall generation quality.

## 6 Conclusion

We introduced Reward-Weighted Sampling (RWS), a novel decoding approach that leverages external reward signals to effectively exploit the non-autoregressive nature of MDMs. Our theoretical analysis demonstrated how reward-based logit scaling can cause beneficial probability rank reversals, fundamentally altering generation order. Empirical evaluations across diverse benchmarks showed that RWS consistently outperforms standard decoding methods, highlighting its ability to enhance generation quality and coherence. Future research may explore adaptive reward scaling, computational efficiency improvements, and extensions to multimodal generation contexts. We hope our work in-

spires further exploration into leveraging reward guidance to fully unlock the potential of diffusion-based language models.

## Limitations

Although Reward-Weighted Sampling (RWS) consistently demonstrates improvements in non-autoregressive generation quality, it introduces additional computational overhead compared to standard confidence-based sampling, as quantitatively analyzed in Appendix N. Specifically, inference time increases by approximately 21–33%, which could impact scalability or real-time applications. Additionally, the performance of RWS inherently depends on the quality and characteristics of the external reward models used. In cases where reward models exhibit biases or misalignment with human preferences, the effectiveness and reliability of RWS could be compromised. Future research should aim to reduce computational overhead and explore robust techniques to select and calibrate reward models effectively.

## Ethical Considerations

Our proposed Reward-Weighted Sampling (RWS) relies on external reward models to guide generation, which could inadvertently propagate biases or subtle misalignments present in those models. Therefore, careful selection and evaluation of reward models are recommended when applying RWS in practice. Future studies might also explore practical methods for monitoring and reducing potential biases in reward-guided text generation systems. For example, simple debiasing measures—such as auditing reward models for group fairness or applying lightweight post-decoding filters—could be incorporated to further mitigate unintended bias (Bu et al., 2025; Allam, 2024).

## References

- Ahmed Allam. 2024. BiasDPO: Mitigating bias in language models through direct preference optimization. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 4: Student Research Workshop)*. Association for Computational Linguistics.
- Jacob Austin, Daniel D. Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. 2021. Structured denoising diffusion models in discrete state-spaces. *CoRR*.
- Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. 2015. Scheduled sampling for sequence prediction with recurrent neural networks. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Yuyan Bu, Liangyu Huo, Yi Jing, and Qing Yang. 2025. Beyond excess and deficiency: Adaptive length bias mitigation in reward models for rlhf. In *Findings of NAACL*.
- Zheng Cai, Maosong Cao, Haojong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, Xiaoyi Dong, Haodong Duan, Qi Fan, Zhaoye Fei, Yang Gao, Jiaye Ge, Chunya Gu, Yuzhe Gu, Tao Gui, and 81 others. 2024. Internlm2 technical report.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2020. Plug and play language models: A simple approach to controlled text generation. In *International Conference on Learning Representations*.
- Prafulla Dhariwal and Alexander Nichol. 2021. Diffusion models beat gans on image synthesis. In *Advances in Neural Information Processing Systems*.
- Nicolai Dorka. 2024. Quantile regression for distributional reward models in rlhf. *arXiv preprint*.
- Marjan Ghazvininejad, Omer Levy, Yinhan Liu, and Luke Zettlemoyer. 2019. Mask-predict: Parallel decoding of conditional masked language models. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*.
- Jiatao Gu, Changhan Wang, and Junbo Zhao. 2019. Levenshtein transformer. In *Advances in Neural Information Processing Systems*.
- Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. 2015. Distilling the knowledge in a neural network. *CoRR*.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In *International Conference on Learning Representations*.

- Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. 2024. Rewardbench: Evaluating reward models for language modeling.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.
- Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori Hashimoto. 2022. Diffusion-LM improves controllable text generation. In *Advances in Neural Information Processing Systems*.
- Bill Yuchen Lin, Wangchunshu Zhou, Ming Shen, Pei Zhou, Chandra Bhagavatula, Yejin Choi, and Xiang Ren. 2020. Commongen: A constrained text generation challenge for generative commonsense reasoning. *Findings of EMNLP*.
- Yen-Ting Lin and Yun-Nung Chen. 2023. LLM-eval: Unified multi-dimensional automatic evaluation for open-domain conversations with large language models. In *Proceedings of the 5th Workshop on NLP for Conversational AI (NLP4ConvAI 2023)*. Association for Computational Linguistics.
- Chris Yuhao Liu, Liang Zeng, Jiacai Liu, Rui Yan, Jujie He, Chaojie Wang, Shuicheng Yan, Yang Liu, and Yahui Zhou. 2024. Skywork-reward: Bag of tricks for reward modeling in llms. *arXiv preprint*.
- Aaron Lou, Chenlin Meng, and Stefano Ermon. 2024. Discrete diffusion modeling by estimating the ratios of the data distribution. In *Proceedings of the 41st International Conference on Machine Learning*.
- Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-Rong Wen, and Chongxuan Li. 2025. Large language diffusion models. *arXiv preprint arXiv:2502.09992*.
- Jingyang Ou, Shen Nie, Kaiwen Xue, Fengqi Zhu, Jiacheng Sun, Zhenguo Li, and Chongxuan Li. 2025. Your absorbing discrete diffusion secretly models the conditional distributions of clean data. In *The Thirteenth International Conference on Learning Representations*.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*.
- Junsoo Park, Seungyeon Jwa, Meiying Ren, Daeyoung Kim, and Sanghyuk Choi. 2024. Offsetbias: Leveraging debiased data for tuning evaluators.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Marc'Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2016. Sequence level training with recurrent neural networks. In *4th International Conference on Learning Representations, ICLR 2016*.
- Alex Wang and Kyunghyun Cho. 2019. BERT has a mouth, and it must speak: BERT as a markov random field language model.
- Kevin Yang and Dan Klein. 2021. FUDGE: Controlled text generation with future discriminators. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- Rui Yang, Ruomeng Ding, Yong Lin, Huan Zhang, and Tong Zhang. 2024. Regularizing hidden states enables learning generalizable reward model for llms. In *Advances in Neural Information Processing Systems*.
- Lifan Yuan, Ganqu Cui, Hanbin Wang, Ning Ding, Xingyao Wang, Jia Deng, Boji Shan, Huimin Chen, Ruobing Xie, Yankai Lin, Zhenghao Liu, Bowen Zhou, Hao Peng, Zhiyuan Liu, and Maosong Sun. 2024. Advancing llm reasoning generalists with preference trees.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging LLM-as-a-judge with MT-bench and chatbot arena. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, and Jiantao Jiao. 2023. Starling-7b: Improving llm helpfulness & harmlessness with rlaif.

## A Masked Diffusion Model Framework

Masked Diffusion Models (MDMs) represent a distinct approach to generative language modeling compared to Autoregressive Models (ARMs). Unlike ARMs, MDMs operate on the entire sequence in a non-autoregressive manner, modeling the data distribution  $p(x_0)$  through a forward masking process and a learned reverse unmasking process (Nie et al., 2025; Ou et al., 2025).

The forward process, denoted  $q(x_t|x_{t-1})$ , gradually corrupts an original clean sequence  $x_0 = (x_0^1, \dots, x_0^L)$  of length  $L$  over a discrete time horizon  $t \in 0, 1, \dots, T$ . Specifically, at step  $t$ , each token  $x_0^j$  in the original sequence is independently replaced by a special mask token,  $M$ , with probability  $t/T$ , or remains unchanged with probability  $1 - t/T$ . This transition probability is defined as:

$$q(x_t^j|x_0^j) = \begin{cases} \frac{t}{T}, & \text{if } x_t^j = M \\ 1 - \frac{t}{T}, & \text{if } x_t^j = x_0^j \end{cases}$$

As  $t$  increases from 0 to  $T$ , the sequence  $x_t$  becomes progressively more masked, until  $x_T$  consists entirely of mask tokens  $M$ .

The reverse process aims to reverse this corruption, generating the clean sequence  $x_0$  starting from the fully masked sequence  $x_T$ . This is achieved by learning a parameterized model, often called a mask predictor  $p_\theta$ , which approximates the conditional probability  $p(x_0|x_t)$ . Typically,  $p_\theta$  is implemented using a non-causal Transformer architecture (Nie et al., 2025). This model takes the masked sequence  $x_t$  at any time  $t$  as input and predicts the original tokens  $x_0$  for all masked positions.

The mask predictor  $p_\theta$  is trained by optimizing an objective derived from the diffusion framework. During training, a time  $t$  is sampled uniformly from  $1, \dots, T$ , the corresponding masked sequence  $x_t$  is generated via the forward process, and the model  $p_\theta$  is trained to predict the original tokens  $x_0$ . The training loss is typically a cross-entropy loss computed only on the masked positions  $j$  where  $x_t^j = M$ .

## B Related Work

**Diffusion Language Models** Diffusion models, initially developed for continuous domains like image generation (Ho et al., 2020; Dhariwal and Nichol, 2021), have recently extended successfully to text. Early diffusion language models (Li et al.,

2022) introduced continuous embedding-based diffusion processes, enabling controllable generation. Subsequently, discrete diffusion approaches (Austin et al., 2021; Lou et al., 2024) emerged, significantly improving text quality. Recent large-scale models like LLaDA (Nie et al., 2025) demonstrated competitive performance against autoregressive models, underscoring the potential of diffusion-based architectures for efficient and flexible text generation.

### Non-Autoregressive Generation Methods

Non-autoregressive (NAR) generation methods, designed to overcome sequential decoding bottlenecks, gained popularity with the introduction of parallel generation techniques like Mask-Predict (Ghazvininejad et al., 2019). Iterative refinement approaches, such as the Levenshtein Transformer (Gu et al., 2019), further enhanced quality and flexibility by iteratively inserting and deleting tokens. These methods established a strong foundation for leveraging iterative masked-decoding strategies, paving the way for diffusion models to exploit non-autoregressive properties effectively.

### Reward-Weighted Decoding Strategies

Recent works integrate external reward signals directly into the decoding process to improve text alignment and quality. Plug-and-play decoding (Dathathri et al., 2020; Yang and Klein, 2021) utilizes auxiliary classifiers to steer autoregressive generation toward desired attributes but faces issues with coherent shifts and error accumulation. Alternatively, reinforcement learning from human feedback (Ouyang et al., 2022) optimizes model parameters using preference-based rewards, though computationally expensive. Direct Preference Optimization (Rafailov et al., 2023) achieves similar goals without RL, fine-tuning models directly with supervised preference signals. Our Reward-Weighted Sampling method extends these ideas to masked diffusion models, incorporating external reward feedback iteratively during decoding process.

## C Reward Normalization Statistics

Table 4 provides the mean ( $\mu_R$ ) and standard deviation ( $\sigma_R$ ) values used for normalizing raw reward scores obtained from each reward model. These statistics were computed from the Nectar validation set (Zhu et al., 2023) and are essential for standardizing reward scales across diverse reward models.

Table 4: Statistics of reward models used for normalization.

Reward Model	Mean	Std
Skywork	-4.95	11.18
LlamaOB	-3.57	2.69
Intern	0.95	1.32
Eurus	954.75	1697.11
GRM	-2.99	3.84
QRM	0.78	0.14

## D Detailed Proof

### D.1 Illustrative Example

Now, we analyze a concrete example to illustrate the theorem. Consider the vectors:

$$\begin{aligned} a &= [1.1, 0.6, 0.3] \\ b &= [1.0, 0.4, 0.4] \end{aligned}$$

First, we compute the gap statistics:

$$\begin{aligned} \Sigma_a &= 1.3, \quad \delta_a = 0.5 \\ \Sigma_b &= 1.2, \quad \delta_b = 0.6 \end{aligned}$$

Since  $\Sigma_a > \Sigma_b$ , no flip is expected for  $r < 1$ . However, because  $\delta_a < \delta_b$ , a flip is expected for  $r > 1$ .

We verify this at specific values:

For  $r = 1$ :

$$P^a(1) \approx 0.486, \quad P^b(1) \approx 0.477$$

Thus,  $P^a(1) > P^b(1)$  as expected.

For  $r = 2$ :

$$P^a(2) \approx 0.835, \quad P^b(2) \approx 0.847$$

Here,  $P^a(2) < P^b(2)$ , confirming a flip for  $r > 1$ .

Solving  $P^a(r) = P^b(r)$  numerically gives  $r^\dagger \approx 1.73$ .

We'll now prove the theorem step by step with detailed explanations.

### D.2 Step 1: Alternative Formulation

First, we'll transform the problem to make it more tractable. Since shifting all logits by a constant doesn't change the softmax probabilities, we can assume without loss of generality that  $a_{t_a} = b_{t_b} = 0$  (i.e., the highest logit in each vector is 0).

Now, let's define:

$$F_a(r) = \sum_{k \neq t_a} e^{-r\Delta_k^a} \quad (4)$$

$$F_b(r) = \sum_{k \neq t_b} e^{-r\Delta_k^b} \quad (5)$$

$$g(r) = F_a(r) - F_b(r) \quad (6)$$

With this notation, the maximum probabilities can be written as:

$$P^a(r) = \frac{1}{1 + F_a(r)} \quad (7)$$

$$P^b(r) = \frac{1}{1 + F_b(r)} \quad (8)$$

Observe that  $P^a(r) > P^b(r)$  if and only if  $F_a(r) < F_b(r)$ , which is equivalent to  $g(r) < 0$ . So:

$$\text{sign}(P^a(r) - P^b(r)) = -\text{sign}(g(r)) \quad (9)$$

This gives us a simpler function  $g(r)$  whose sign tells us which vector has the higher maximum probability.

### D.3 Step 2: Behavior for Small $r$

We want to understand how  $g(r)$  behaves as  $r$  approaches 0. For small values of  $r$ , we can use the Taylor series expansion of the exponential function:

$$e^{-r\Delta} \approx 1 - r\Delta + \frac{1}{2}r^2\Delta^2 + O(r^3) \quad (10)$$

Applying this to our functions  $F_a(r)$  and  $F_b(r)$ :

$$\begin{aligned} F_a(r) &\approx \sum_{k \neq t_a} (1 - r\Delta_k^a) \\ &\quad + \frac{1}{2}r^2(\Delta_k^a)^2 + O(r^3) \Big) \\ &= (d-1) - r \sum_{k \neq t_a} \Delta_k^a \\ &\quad + \frac{1}{2}r^2 \sum_{k \neq t_a} (\Delta_k^a)^2 + O(r^3) \\ &= (d-1) - r\Sigma_a + O(r^2) \end{aligned} \quad (11)$$

Similarly:

$$F_b(r) \approx (d-1) - r\Sigma_b + O(r^2) \quad (12)$$

Therefore:

$$\begin{aligned} g(r) &= F_a(r) - F_b(r) \\ &\approx ((d-1) - r\Sigma_a) \\ &\quad - ((d-1) - r\Sigma_b) + O(r^2) \\ &= r(\Sigma_b - \Sigma_a) + O(r^2) \end{aligned} \quad (13)$$

This tells us that for small positive  $r$ :

$$\text{sign}(g(r)) \approx \text{sign}(\Sigma_b - \Sigma_a) \quad (14)$$

So if  $\Sigma_a < \Sigma_b$ , then  $g(r) > 0$  for small  $r$ , which means  $P^a(r) < P^b(r)$ . But we know that at  $r = 1$ , we have  $P^a(1) > P^b(1)$  (our initial assumption), which means  $g(1) < 0$ .

Since  $g(r)$  is continuous (as it's a difference of sums of continuous functions), if  $g(r) > 0$  for small  $r$  and  $g(1) < 0$ , then by the Intermediate Value Theorem, there must exist some  $r^* \in (0, 1)$  where  $g(r^*) = 0$ , which means  $P^a(r^*) = P^b(r^*)$ .

Moreover, since  $g(r) > 0$  for small positive  $r$ , we have  $P^a(r) < P^b(r)$  for all  $0 < r < r^*$ .

#### D.4 Step 3: Behavior for Large $r$

For large values of  $r$ , the exponential terms with the smallest exponents will dominate the sums. Let's analyze this behavior:

$$F_a(r) = \sum_{k \neq t_a} e^{-r\Delta_k^a} \quad (15)$$

$$= e^{-r\delta_a} \cdot \left( 1 + \sum_{k \neq t_a, \Delta_k^a > \delta_a} e^{-r(\Delta_k^a - \delta_a)} \right) \quad (16)$$

As  $r \rightarrow \infty$ , all the terms in the sum approach 0, except for the terms where  $\Delta_k^a = \delta_a$ . Let's say there are  $C_a$  such terms. Then:

$$F_a(r) \sim C_a \cdot e^{-r\delta_a} \quad \text{as } r \rightarrow \infty \quad (17)$$

Similarly:

$$F_b(r) \sim C_b \cdot e^{-r\delta_b} \quad \text{as } r \rightarrow \infty \quad (18)$$

Therefore:

$$g(r) = F_a(r) - F_b(r) \quad (19)$$

$$\sim C_a \cdot e^{-r\delta_a} - C_b \cdot e^{-r\delta_b} \quad \text{as } r \rightarrow \infty \quad (20)$$

If  $\delta_a > \delta_b$ , then as  $r \rightarrow \infty$ , the term  $C_b \cdot e^{-r\delta_b}$  dominates, and:

$$g(r) \sim -C_b \cdot e^{-r\delta_b} < 0 \quad \text{as } r \rightarrow \infty \quad (21)$$

This means that for sufficiently large  $r$ , we have  $g(r) < 0$ , which implies  $P^a(r) > P^b(r)$ .

But if  $\delta_a < \delta_b$ , then as  $r \rightarrow \infty$ , the term  $C_a \cdot e^{-r\delta_a}$  dominates, and:

$$g(r) \sim C_a \cdot e^{-r\delta_a} > 0 \quad \text{as } r \rightarrow \infty \quad (22)$$

So if  $\delta_a < \delta_b$ , then for sufficiently large  $r$ , we have  $g(r) > 0$ , which implies  $P^a(r) < P^b(r)$ .

Now, we know that at  $r = 1$ , we have  $P^a(1) > P^b(1)$  (our initial assumption), which means  $g(1) < 0$ .

If  $\delta_a > \delta_b$ , then for sufficiently large  $r$ , we have  $g(r) < 0$ , so  $P^a(r) > P^b(r)$ . There's no flip for large  $r$ .

But if  $\delta_a < \delta_b$ , then for sufficiently large  $r$ , we have  $g(r) > 0$ , so  $P^a(r) < P^b(r)$ . Since  $g(1) < 0$  and  $g(r) > 0$  for large  $r$ , by the Intermediate Value Theorem, there must exist some  $r^\dagger > 1$  where  $g(r^\dagger) = 0$ , which means  $P^a(r^\dagger) = P^b(r^\dagger)$ .

#### D.5 Step 4: Uniqueness and Exclusivity of the Cases

We need to show that the scenarios described in the theorem are mutually exclusive and that the crossing points  $r^*$  and  $r^\dagger$  are unique.

First, note that  $P^a(r)$  and  $P^b(r)$  are both strictly monotonic functions of  $r$ . Specifically, as  $r$  increases, the maximum probabilities also increase (the distributions become more concentrated on the highest logits).

Since  $P^a(r)$  and  $P^b(r)$  are both strictly increasing functions of  $r$ , their difference  $P^a(r) - P^b(r)$  can change sign at most once as  $r$  increases from 0 to  $\infty$ . This means that the function  $g(r)$  can also change sign at most once.

So if  $g(r)$  changes from positive to negative as  $r$  increases (which happens when  $\Sigma_a < \Sigma_b$ ), then it cannot change back to positive. Similarly, if  $g(r)$  changes from negative to positive (which happens when  $\delta_a < \delta_b$ ), then it cannot change back to negative.

This means that the two scenarios (low- $r$  flip and high- $r$  flip) are mutually exclusive.

Moreover, if neither gap condition holds (i.e., if  $\Sigma_a \geq \Sigma_b$  and  $\delta_a \geq \delta_b$ ), then  $g(r)$  cannot change sign. Since  $g(1) < 0$  (our initial assumption), this means  $g(r) < 0$  for all  $r > 0$ , which implies  $P^a(r) > P^b(r)$  for all  $r > 0$ . This is the "no flip" scenario.

## E Proof of Theorem 2

Here we restate and provide a detailed proof of Theorem 2 from the main text.

**Theorem 3** (Reward-Monotonic Guidance). *Let  $p_\theta(x) \propto \exp L_\theta(x)$  be the base decoder and  $p_r(x) \propto p_\theta(x) e^{rR(x)}$  its reward-scaled variant*

with any  $r > 0$ . Then

$$\mathbb{E}_{p_r}[R] \geq \mathbb{E}_{p_\theta}[R], \quad (23)$$

with strict inequality if  $\text{Var}_{p_r}[R] > 0$ .

*Proof.* Define the log-partition function as:

$$A(r) = \log Z_r = \log \sum_x e^{L_\theta(x) + r R(x)}.$$

Using standard exponential-family calculus, we have:

$$A'(r) = \sum_x p_r(x) R(x) = \mathbb{E}_{p_r}[R],$$

and the second derivative is:

$$A''(r) = \text{Var}_{p_r}[R] \geq 0.$$

Since  $A''(r) \geq 0$ , the function  $A'(r)$  is non-decreasing with respect to  $r$ . Integrating  $A''$  from 0 to  $r$ , we obtain:

$$\begin{aligned} A'(r) &= A'(0) + \int_0^r A''(s) ds \\ &= \mathbb{E}_{p_\theta}[R] + r \text{Var}_{p_r}[R]. \end{aligned}$$

This directly implies the desired result:

$$\mathbb{E}_{p_r}[R] = \mathbb{E}_{p_\theta}[R] + r \text{Var}_{p_r}[R] \geq \mathbb{E}_{p_\theta}[R],$$

with strict inequality whenever  $\text{Var}_{p_r}[R] > 0$ .  $\square$

## F Comparison with Softmax Temperature Scaling

Reward-Weighted Sampling (RWS) leverages global feedback from an external reward model to adaptively scale logits at each diffusion step. A related but distinct technique is softmax temperature scaling, which uniformly scales logits using a fixed temperature parameter across the entire decoding process. Here, we highlight both theoretical and empirical differences between RWS and fixed temperature scaling.

The primary theoretical distinction is adaptivity. Temperature scaling applies a constant scaling factor, irrespective of intermediate sequence quality or global coherence signals, potentially causing overly aggressive token selections or negligible impact. In contrast, RWS dynamically adjusts scaling based on global reward evaluations, ensuring that logits are adjusted meaningfully according to generation quality.

Empirically, we conducted a comparative experiment using the keyword-constrained generation task (Section 5.4) with the same hyperparameter search space  $\{0.01, 0.1, 1, 2, 4, 8, 16, 32\}$  for temperature scaling. Perplexity scores obtained were consistently higher (worse) than those achieved by RWS across all reward models. Specifically, temperature scaling resulted in perplexities ranging from 11.03 to 14.76, all worse than the worst-performing RWS configuration (10.84). Interestingly, Generation Order Deviation (GOD) values for temperature scaling ranged from 0.79 to 2.31, overlapping substantially with RWS’s performance range. This demonstrates that while temperature scaling can similarly induce non-autoregressive behavior, it does so at the significant expense of generation quality.

In summary, adaptive reward-based scaling employed by RWS provides substantial qualitative improvements over simple, fixed temperature scaling, demonstrating the critical importance of incorporating meaningful global reward signals into the decoding process (see Appendix G for all detailed results).

## G Full Results of Softmax Temperature Scaling

We provide detailed results comparing different softmax temperature scaling values on the keyword-constrained generation task (Section 5.4). Table 5 summarizes perplexity and Generation Order Deviation (GOD) values across all tested temperature hyperparameters.

Temperature	Perplexity ( $\downarrow$ )	GOD ( $\uparrow$ )
0.01	13.57	0.91
0.1	13.52	0.95
1.0	14.76	0.80
2.0	11.72	1.56
4.0	11.13	2.04
8.0	11.13	2.26
16.0	11.03	<b>2.31</b>
32.0	11.58	1.88

Table 5: Perplexity and Generation Order Deviation (GOD) of softmax temperature scaling on the keyword-constrained generation task. While higher temperature scaling values increase non-autoregressive characteristics (higher GOD), they consistently degrade fluency (higher perplexity) compared to Reward-Weighted Sampling (RWS).

Name	Model	Size	Base
Skywork	Skywork-Reward-Llama-3.1	8B	LLaMA-3.1
LlamaOB	Llama-3-OffsetBias-RM	8B	LLaMA-3
Intern	InternLM2-Reward	7B	InternLM
Eurus	Eurus-RM	7B	LLaMA-2
GRM	GRM-Llama3.2	3B	LLaMA-3.2
QRM	QRM-Llama3.1	8B	LLaMA-3.1

Table 6: Characteristics of reward models used for evaluation and guidance

## H Details of Reward Models

Table 6 summarizes the reward models used for both guidance and evaluation throughout our experiments. The models are selected from RewardBench (Lambert et al., 2024) and vary in size, base architecture, and training methodologies.

## I Detailed Experimental Setup

**Diffusion LLM.** We use LLaDA-8B-Instruct as our base Masked Diffusion Model (MDM) for all experiments. This non-autoregressive model employs a diffusion-based architecture that iteratively unmasks tokens, as described in Appendix A.

**Reward Models.** We select six publicly available reward models (Liu et al., 2024; Park et al., 2024; Cai et al., 2024; Yuan et al., 2024; Dorka, 2024; Yang et al., 2024) of varying sizes and architectures (see Table 6) from RewardBench (Lambert et al., 2024).<sup>3</sup> We choose models balancing performance and computational feasibility given our resources. Each reward model was previously trained to predict human preferences across diverse text generation tasks, and their outputs are normalized using pre-computed statistics shown in Table 4.

**Datasets.** We evaluate our approach on three benchmarks:

1. **RewardBench** (Lambert et al., 2024): comprehensive evaluation of helpfulness, harmlessness, and honesty.
2. **MT-Bench** (Zheng et al., 2023): multi-turn instruction-following across diverse domains.
3. **Keyword-Constrained Generation Task**: specifically designed for this study, requiring incorporation of specific keywords while maintaining coherence, testing flexible token selection (dataset construction detailed in Appendix M).

<sup>3</sup><https://huggingface.co/spaces/allenai/reward-bench>

**Additional Metrics.** We also employ additional standard metrics alongside GOD:

- *Win Rate*: Proportion of prompts where a method’s output is rated higher by reward models or LLM judges.
- *Perplexity*: Standard measure of text fluency and coherence, particularly relevant in keyword-constrained tasks.

**General Implementation Details.** Unless otherwise specified, general experimental settings are:

- Diffusion steps: 128
- Block size: 32 tokens
- Maximum output length: 256 tokens
- Optimal reward scales individually selected per reward model from the set {0.01, 0.1, 1.0, 2.0, 4.0, 8.0, 16.0, 32.0} using the RewardBench validation set.
- Hardware: NVIDIA A100 GPUs (40GB VRAM)

Reward signals were normalized using statistics summarized in Table 4.

## J LLM-as-a-judge Evaluation Details

For multi-turn dialogue evaluations using MT-Bench, we employ GPT-4.1 (OpenAI API) as a judge model to determine which of two provided responses is of higher quality. The exact prompt used for the evaluation is presented in Figure 7.

<b>System Prompt:</b> You are a strict evaluator. Decide which answer (1 or 2) is higher quality. If answer 1 is better, respond with 1. If answer 2 is better, respond with 2. Respond ONLY with one character: 1 or 2.
<b>User Prompt:</b> Answer 1: <i>[First model-generated response here]</i> Answer 2: <i>[Second model-generated response here]</i>

Figure 7: Exact prompt template used for GPT-4.1 judge evaluations. The model compares two responses and selects the higher-quality response.

For each comparison, GPT-4.1 returns either "1" or "2" directly indicating the preferred response.

## K Detailed LLM Judge Results

Table 7 presents detailed results of the LLM-as-a-judge evaluations, showing win rates separately for each ordering configuration.

## L Qualitative Multi-turn Examples

In this section, we present qualitative examples contrasting Reward-Weighted Sampling (RWS) and Default Sampling methods in multi-turn dialogue settings.

**Brief Comparison** Figure 8 provides a concise comparison highlighting the key differences in coherence and detail between the two methods across a two-turn dialogue scenario.

**Detailed Comparison (Full Version)** For a more comprehensive illustration, Figure 10 shows the complete dialogue interactions, clearly demonstrating how RWS effectively maintains coherence and context across multiple turns compared to Default Sampling.

## M Dataset Procedure

To construct the keyword-constrained generation dataset used in our experiments, we adopted a systematic procedure to ensure topic diversity, task relevance, and quality:

**Step 1: Topic Selection.** We selected 10 representative topics across various domains to ensure comprehensive coverage of real-world scenarios and broad generalization. The selected topics were:

1. Science
2. Technology
3. Economy
4. Health
5. Sports
6. Environment
7. Politics
8. Education
9. Culture
10. Entertainment

**Step 2: Prompt and Keyword Generation.** Using the GPT-4.1 API, we generated an initial pool of 100 keyword-constrained prompts, equally distributed across the 10 selected topics (10 prompts per topic). Each prompt was designed to incorporate specific anchor keywords naturally within coherent and contextually appropriate responses.

## Step 3: Quality Control and Expert Validation.

The generated dataset underwent rigorous review by three PhD students. They independently evaluated each prompt based on difficulty and overall quality. Prompts identified with issues such as duplication, overly simplistic difficulty, grammatical errors, ambiguity, or poor contextual quality were flagged for removal.

**Step 4: Refinement and Final Selection.** Based on expert validation, 15 prompts were excluded due to quality concerns. Subsequently, an additional set of prompts was generated using GPT-4.1, maintaining the same distribution across topics. After this iterative refinement, we finalized a dataset consisting of 100 keyword-constrained prompts, equally distributed among the 10 selected topics.

This structured methodology ensured the dataset was both diverse and representative, suitable for effectively evaluating the non-autoregressive generation capabilities of our model.

## N Computational Overhead Analysis

To evaluate the additional computational overhead introduced by our proposed Reward-Weighted Sampling (RWS) method compared to the standard confidence-based sampling baseline, we conducted experiments measuring inference time and GPU memory usage.

Experiments were performed on a single NVIDIA A100 GPU (40GB VRAM) with batch size 1, using the keyword-constrained generation dataset (Experiment 3). We measured the average inference time per sample across 200 prompts randomly selected from this dataset. GPU memory usage was recorded using PyTorch’s built-in memory monitoring functions, with additional GPU memory usage primarily attributed to the use of reward models.

Results are summarized in Table 8. We observed an inference time increase ranging from approximately 21% to 33% with RWS compared to the baseline method, varying by the specific reward model used. GPU memory overhead due to reward models varied significantly, showing an increase ranging from approximately 6.5 GiB to 14.8 GiB depending on the reward model.

These results indicate that RWS introduces a manageable computational overhead, maintaining practicality for real-world applications while significantly enhancing generation quality.

Reward Model	RWS first (%)		RWS second (%)	
	RWS Win	Default Win	RWS Win	Default Win
Skywork	<b>65.0</b>	35.0	<b>65.0</b>	35.0
LlamaOB	<b>62.5</b>	37.5	<b>62.5</b>	37.5
Intern	<b>61.25</b>	38.75	<b>60.0</b>	40.0
Eurus	<b>65.0</b>	35.0	<b>58.75</b>	41.25
GRM	<b>65.0</b>	35.0	<b>61.25</b>	38.75
QRM	<b>60.0</b>	40.0	<b>60.0</b>	40.0

Table 7: Full LLM judge-based evaluation results for each response ordering configuration.

Method	Avg. Inference Time (s)	Additional GPU Memory (MiB)
Default	20.197	-
RWS (Skywork)	26.288	14,812
RWS (LlamaOB)	26.046	14,812
RWS (Intern)	26.560	14,502
RWS (Eurus)	26.950	14,078
RWS (GRM)	24.597	6,542
RWS (QRM)	26.330	14,812

Table 8: Computational overhead comparison between confidence-based sampling and Reward-Weighted Sampling (RWS).

## O Artifact and Data Usage

**License for Artifacts.** All reward models and datasets employed in this research, including RewardBench (Lambert et al., 2024) and MT-Bench (Zheng et al., 2023), are publicly available under permissive licenses for research purposes. Specifically, RewardBench is released under an MIT License, and MT-Bench is provided under Apache 2.0 License.

**Artifact Use Consistent with Intended Use.** The artifacts used, including publicly available reward models from RewardBench, are employed strictly in alignment with their intended research-oriented use as clearly specified in their original publications. No derivative artifacts were created or distributed beyond the research context, complying with the original access conditions.

**Personally Identifying Information and Offensive Content.** The datasets and prompts utilized in our experiments (RewardBench, MT-Bench, and keyword-constrained generation tasks) were reviewed to ensure they do not contain personally identifying information or offensive content. All prompts underwent careful expert validation and refinement (Appendix M), explicitly removing any content that could uniquely identify individuals or

be considered offensive.

**Documentation of Artifacts.** Comprehensive documentation for each artifact (RewardBench and MT-Bench) is provided by their original authors, clearly detailing domains, languages, and linguistic phenomena covered. Our custom-created keyword-constrained generation dataset includes clear documentation on domains, representative topics, and the validation process, as described in Appendix M.

## P Use of AI Assistants

We used ChatGPT to help refine manuscript wording, improve clarity, and assist with LaTeX formatting. All research, experiments, and analyses are solely the authors' work.

Turn	Default Sampling (Baseline)	Reward-Weighted Sampling (Ours)
Prompt 1	<b>Explain what's base rate fallacy and list five specific examples of how politicians use it for campaigns.</b>	
Answer 1	<ul style="list-style-type: none"> <li>- Defines base rate fallacy adequately but provides repetitive examples: "Highlighting candidate's past successes to suggest future success." (Repeated)</li> <li>"Using candidate's past failures to suggest future failures." (Repeated)</li> <li>- Includes logically inconsistent example: "Highlighting opponent's past successes."</li> </ul>	<ul style="list-style-type: none"> <li>- Clearly defines base rate fallacy: "Overestimating or underestimating base rates of events."</li> <li>- Provides diverse examples: <ul style="list-style-type: none"> <li>1. "Highlighting negative instances selectively."</li> <li>2. "Exaggerating own policy successes."</li> <li>3. "Presenting misleading statistics."</li> <li>4. "Ignoring evidence selectively."</li> <li>5. "Creating false biases favoring own policies."</li> </ul> </li> </ul>
Prompt 2	<b>Provide a detailed plan for an election campaign using the first example.</b>	
Answer 2	<ul style="list-style-type: none"> <li>- Repetitive, vague steps without actionable details: "Emphasize the candidate's past successes repeatedly."</li> <li>- Lacks concrete strategies.</li> </ul>	<ul style="list-style-type: none"> <li>- Coherent and actionable strategies clearly aligned with the example: <ul style="list-style-type: none"> <li>1. "Highlight opponent's negative past events through media/social media."</li> <li>2. "Downplay own positive records selectively."</li> <li>3. "Use emotional appeals in negative campaigns."</li> <li>4. "Leverage media to amplify opponent's negatives."</li> <li>5. "Conduct targeted social media outreach."</li> </ul> </li> </ul>

Figure 8: Comparison of multi-turn responses between Default Sampling (Baseline) and Reward-Weighted Sampling (Ours). Our method produces clearer definitions, diverse examples, and coherent follow-up plans, showcasing the effectiveness of reward guidance in improving logical consistency and practical applicability. Full responses are provided in Figure 10.

Turn	Default Sampling (Baseline)	Reward-Weighted Sampling (Ours)
<b>Prompt 1</b>	Explain what's base rate fallacy and list five specific examples of how politicians use it for campaigns.	
<b>Answer 1</b>	<p>The base rate fallacy is a cognitive bias that occurs when an individual judges the probability of an event based on the frequency of similar event, rather than the actual probability of the event occurring. Five examples of how politicians use the base rate fallacy for campaigns are:</p> <ol style="list-style-type: none"> <li>1. Highlighting to the success of a candidate's past policies to suggest that they current policy will be successful.</li> <li>2. Using the candidate's past failures to win promises to suggest that they are not capable of keeping future promises.</li> <li>3. Emphasizing the candidate's past successes to suggest that they will continue to be successful in future endeavors.</li> <li>4. Using the candidate's past failures to suggest that they will fail in future future endeavors.</li> <li>5. Highlighting the opposing candidate's past successes to suggest that they will be successful in the future.</li> </ol>	<p>The base rate fallacy is a cognitive bias where people overestimate the base rate of an event or underestimate the base rate of another event. Five examples of how politicians use the base rate fallacy for campaigns are:</p> <ol style="list-style-type: none"> <li>1. <b>Highlighting negative instances:</b> Politicians may use the base rate fallacy by highlighting negative instances of their opponents while ignoring the positive instances of their own policies.</li> <li>2. <b>Exaggerating successes:</b> Politicians may use the base base rate fallacy by exaggerating the successes of their own policies while downplaying the failures of their opponents' policies.</li> <li>3. <b>Misleading statistics:</b> Politicians may use the base rate fallacy by presenting misleading statistics or statistics that are manipulated to support their argument.</li> <li>4. <b>Ignoring the evidence:</b> Politicians may use the base base base rate fallacy by ignoring the evidence against their policies while focusing on positive evidence.</li> <li>5. <b>False bias:</b> Politicians may use the base base base base rate fallacy by creating a false bias in favor of their own policies while ignoring the evidence against their policies.</li> </ol>

Figure 9: Complete responses to Prompt 1 from both Default Sampling and Reward-Weighted Sampling methods, as referenced in Figure 8.

Turn	Default Sampling (Baseline)	Reward-Weighted Sampling (Ours)
Prompt 2	Provide a detailed plan for an election campaign using the first example.	
Answer 2	<p>1. Develop a well and campaign message that emphasizes the candidate's past successes and how they will continue to be successful in the future. 2. Create a campaign website that highlights the candidate's past successes and how they will continue to be successful in the future. 3. Launch a media campaign that emphasizes the candidate's past successes and how they will continue to be successful in the future. 4. Hold rallies and rallies that emphasize the candidate's past successes and how they will continue to be successful in the future. 5. Organize and events that emphasize the candidate's past successes and how they will continue to be successful in the future. 6. Launch a fundraising campaign that emphasizes the candidate's past successes and how they will continue to be successful in the future. 7. Create a website that emphasizes the candidate's past successes successes and how they will continue to be successful in the future. 8. Launch a media campaign that emphasizes the candidate's past successes and how they will continue to be successful in the future. 9. Hold public events that emphasize the candidate's past successes and how they will continue to be successful in the future. 10. Promote the campaign by launching a social media campaign that emphasizes the candidate's past successes.</p>	<p>A detailed plan for an election campaign campaign using the first example of the base rate fallacy could include:</p> <ol style="list-style-type: none"> <li>1. Highlighting negative instances of the opponent: The campaign team could highlight negative instances instances of the opponent's policies, such as past failures, missteps, or mis mistakes. This could be done through media coverage, social media, and public events.</li> <li>2. Ignoring positive instances of the candidate's policies: The campaign team could ignore positive instances of the candidate's policies, such as past successes or positive outcomes. This could be done by not mentioning them in the campaign materials or by downplaying their significance.</li> <li>3. Using emotional appeals: The campaign team could use emotional appeals to highlight the negative instances of the opponent's policies. This could be done by using stories or statistics that highlight the negative impact of the opponent's policies.</li> <li>4. Using media coverage: The campaign team could use media coverage to highlight the negative instances of the opponent's policies. This could be done by reaching out to media outlets and providing them with negative instances of the opponent's policies.</li> <li>5. Using social media: The campaign team could use social media to highlight the negative instances of the opponent's policies. This could be done by creating posts.</li> </ol>

Figure 10: Complete responses to Prompt 2 from both Default Sampling and Reward-Weighted Sampling methods, as referenced in Figure 8.