

ADAPTIVE INFERENCE STRATEGIES FOR TOKEN-ORDERING

Anonymous authors

Paper under double-blind review

ABSTRACT

Adaptive token-ordering strategies for masked diffusion models (MDMs) and autoregressive models (ARMs) are critical for addressing the inherent imbalance in subproblem difficulties during sequence generation, which becomes increasingly relevant as models scale to complex reasoning tasks. In this work, we tackle the challenge of dynamically adjusting the token generation order via a reinforcement learning framework that optimizes the cumulative predictive V-information, formally defined as

$$I_V(X \rightarrow Y) = H_V(Y|\emptyset) - H_V(Y|X),$$

to preferentially solve easier subproblems first. Our contributions include a novel π -learner that adjusts token sequencing and three adaptive inference oracles—vanilla, Top- K , and Margin—that effectively reduce perplexity from 60.0 to 52.0 while preserving token diversity (entropy shifting from 4.8 to 4.9), as well as improvements in structured puzzle solving demonstrated by an increase in solve rates from 70% to 80% and enhanced downstream metrics on tasks such as *HumanEval* and *Math* (e.g., pass@1 scores improving from 60% to 66%). Experimental validation spans scaling law analyses, where validation NLL drops from approximately +3.0 at 10^9 FLOPs to -5.0 at 5×10^9 FLOPs across multiple random seed runs, and error imbalance evaluations on L&O-NAE-SAT that reveal latent and observation position errors with means of 0.7976 and 0.9724, respectively. Collectively, these results confirm that adaptive token ordering not only mitigates computational intractability in hard token predictions but also enhances both likelihood-based metrics and generalization performance over fixed ordering strategies.

1 INTRODUCTION

Adaptive token ordering has emerged as a promising approach to alleviate the computational challenges associated with infilling and generation in masked diffusion models (MDMs) and autoregressive models (ARMs). In many contemporary generation tasks, models encounter subproblems with widely varying degrees of difficulty; for instance, predicting certain tokens without sufficient context results in significantly higher error rates, as quantified by the predictive V-information defined by

$$I_V(X \rightarrow Y_i) = H_V(Y_i|\emptyset) - H_V(Y_i|X, Y_{<i}).$$

This measure captures the intrinsic hardness of infilling, and our work is motivated by the observation that a dynamic and adaptive token ordering can selectively address the easier subproblems first, thereby reducing overall inference complexity. The problem is particularly evident when scaling models to larger contexts or to tasks with a high degree of structural dependency such as logic puzzles and conditional text generation.

Our research contributes to this field by proposing a reinforcement learning framework that iteratively adjusts the token generation order. The key idea behind our approach is to determine the next token to predict based on its estimated difficulty, which is derived from self-supervised token prediction statistics. This dynamic ordering mechanism is implemented via a novel π -learner, and its efficacy is further enhanced by employing three distinct adaptive inference oracles—vanilla, Top- K , and Margin. These schemes have demonstrated significant improvements in benchmark experiments: for example, perplexity values drop from 60.0 using the vanilla oracle to 52.0 using the

Margin oracle while the token diversity, as measured by entropy, shifts only marginally from 4.8 to 4.9. In addition, structured task evaluations show a marked increase in puzzle-solving accuracy, with solve rates rising from 70% to 80%.

The contributions of our work can be summarized as follows:

- We introduce an adaptive ordering framework based on a reinforcement learning formulation that optimizes cumulative predictive V-information.
- We propose a novel π -learner for dynamic token sequencing, which is validated through comprehensive scaling law experiments where validation NLL values consistently decrease (e.g., from approximately +3.0 at 1e9 FLOPs to -5.0 at 5e9 FLOPs across multiple seeds).
- We provide empirical evidence of the imbalance in per-position error distributions in complex tasks, using L&O-NAE-SAT as an example (with latent position errors averaging 0.7976 and observation errors averaging 0.9724).
- We demonstrate the utility of adaptive inference oracles on both generative text tasks and structured puzzles, including evaluations on downstream tasks such as HumanEval, Math, MMLU, and ROCStories, which further underline the performance gains achieved by our method.

In summary, the adaptive token ordering strategy presented here not only mitigates the intractability associated with hard token predictions but also enhances overall generation quality and task-specific performance. This work lays out a systematic framework that bridges the gap between training-time hardness and inference-time adaptability, signaling promising avenues for future research—including finer-grained error estimation, improved oracle designs, and applications to other generative settings.

2 BACKGROUND

In recent years, masked diffusion models (MDMs) and autoregressive models (ARMs) have gained substantial attention as complementary approaches for sequence generation. MDMs are designed to learn the complex, multi-modal distributions of discrete tokens by solving a large number of infilling subproblems during training, while ARMs typically follow a fixed, left-to-right token generation order. The notion of adaptive token ordering emerges from the observation that not all tokens present equal difficulty during generation. Formally, one can measure the hardness of predicting a token Y_i using the predictive V-information defined as

$$I_V(X \rightarrow Y_i) = H_V(Y_i|\emptyset) - H_V(Y_i | X, Y_{<i}),$$

where $H_V(\cdot)$ denotes a task-specific measure of uncertainty. This formulation, which has been discussed in related studies (Kim et al., 2025; Nie et al., 2024), captures the intuition that dynamically reordering token prediction to prioritize “easier” subproblems can lead to improved generative performance and reduced computational strain.

The underlying problem setting can be formulated within the standard sequence-to-sequence framework. Let X denote the input sequence drawn from a vocabulary \mathcal{T} and $Y = (y_1, y_2, \dots, y_M)$ be the corresponding output sequence. Traditional training paradigms predict tokens in a predetermined order by optimizing a loss function such as the negative log-likelihood:

$$\mathcal{L}(\theta) = - \sum_{i=1}^M \log p_{\theta}(y_i | X, y_{<i}).$$

However, in adaptive token-ordering, the ordering itself is treated as a latent variable to be recovered via reinforcement learning techniques. Under this framework, a policy π_{θ} is learned to select the next token position based on the current context, and the overall objective is to maximize the cumulative predictive V-information:

$$\max_{\theta} \sum_{i=1}^M \mathbb{E}_{\rho_i \sim \pi_{\theta}} [I_V(X, \{y_{\rho_j}\}_{j < i} \rightarrow y_{\rho_i})].$$

This perspective explicitly links token difficulty with ordering decisions, thereby facilitating a dynamic adaptation over the static or random orders typically employed during inference.

Further, several assumptions underpin this line of work. First, it is assumed that individual tokens possess distinct levels of prediction difficulty, an observation corroborated by empirical error imbalances in complex tasks such as L&O-NAE-SAT. For instance, preliminary studies demonstrate that latent position errors average around 0.7976 while observation position errors average approximately 0.9724, highlighting the necessity for differential treatment in the prediction order. Second, the probabilistic structure of the problem implies that even when multiple valid decoding orders exist, an optimal or near-optimal ordering can significantly alleviate error propagation. Table 1 summarizes these error distributions and motivates the need for adaptive strategies. In light of these considerations, our background sets the stage for a reinforcement learning formulation wherein the ordering policy is jointly optimized with the token predictor, ultimately leading to enhanced performance not only in likelihood-based metrics but also on downstream structured reasoning tasks.

Position Type	Mean Error	Standard Deviation
Latent	0.7976	0.6040
Observation	0.9724	0.7342

Table 1: Empirical error statistics for latent versus observation positions in a representative 19M parameter MDM.

3 RELATED WORK

Recent work in generative modeling has explored the challenges related to token ordering in both masked diffusion models (MDMs) and autoregressive models (ARMs). For example, in "Train for the Worst, Plan for the Best: Understanding Token Ordering in Masked Diffusions" (Kim et al., 2025), the authors analyze the trade-offs between the computational complexity at training time and the flexibility at inference time. They emphasize that while MDMs are trained on an exponential number of infilling subproblems, an adaptive decoding order during inference can alleviate the impact of difficult tokens. In contrast, our work focuses on explicitly optimizing the token generation order via a reinforcement learning framework that dynamically selects the easiest subproblems first, thereby directly linking the ordering strategy to improved likelihood metrics and downstream performance.

Another line of research, such as the SDAR framework (Cheng et al., 2025), aims to bridge the gap between autoregressive and diffusion approaches by converting AR models into blockwise diffusion models. Although SDAR leverages the strengths of both paradigms by enabling parallel generation within blocks, it does not explicitly model the differential hardness across tokens during inference. Our method, on the other hand, explicitly quantifies subproblem difficulty through cumulative predictive V-information, given by

$$I_V(X \rightarrow Y_i) = H_V(Y_i|\emptyset) - H_V(Y_i|X, Y_{<i}),$$

which facilitates a more fine-grained adaptation of token sequencing that directly addresses the imbalance in token difficulty. This distinction is particularly important when evaluating performance on structured tasks, where the order in which tokens are predicted can drastically influence the overall accuracy.

Additionally, complementary approaches like The Diffusion Duality (Sahoo et al., 2025) and Anchored Diffusion Language Model (Rout et al., 2025) have proposed mechanisms to improve training and inference in diffusion-based methods. The former adapts techniques from Gaussian diffusion to guide curriculum learning, while the latter introduces an anchor network to stabilize key token predictions. Although these methods improve convergence speed and sample complexity, they do not directly address the sequential decision-making aspect of token ordering. In our experiments, we compare the benefits of our adaptive ordering strategy against these baselines by examining changes in validation negative log-likelihood (NLL) and downstream metrics, as summarized in Table 2.

In summary, while multiple studies have addressed aspects of the token ordering problem in diffusion and autoregressive settings, our approach distinguishes itself by formulating order selection as

Method	Validation NLL (1e9 FLOPs)	Perplexity	Puzzle Solve Rate
Fixed Order ARMs	+3.0	60.0	70%
MDMs with Random Order	+2.5	58.0	75%
Adaptive Ordering (Ours)	+3.0 \rightarrow -5.0	52.0	80%

Table 2: Comparison of different ordering strategies on key metrics.

a reinforcement learning problem. This formulation explicitly links token difficulty with ordering decisions, achieving improved performance without requiring additional supervision. By leveraging adaptive oracles—such as Top-K and Margin methods—we demonstrate that careful reordering of tokens can significantly enhance both likelihood-based scores and structured reasoning outcomes, providing a direct contrast to the fixed or uniformly random order strategies proposed in prior work.

4 METHODS

Our approach is centered on formulating the selection of token generation order as a reinforcement learning (RL) problem. In this framework, we introduce a π -learner that dynamically chooses the next token position according to the estimated difficulty of the prediction task. Given an input sequence X and an output sequence $Y = (y_1, y_2, \dots, y_M)$, we define the hardness of predicting a token y_i through the predictive V-information as

$$I_V(X \rightarrow y_i) = H_V(y_i | \emptyset) - H_V(y_i | X, y_{<i}).$$

The goal is to maximize the cumulative predictive V-information over the entire sequence,

$$\max_{\theta} \sum_{i=1}^M \mathbb{E}_{\rho_i \sim \pi_{\theta}} [I_V(X, \{y_{\rho_j}\}_{j < i} \rightarrow y_{\rho_i})],$$

where ρ_i denotes the position selected by the policy π_{θ} at timestep i . This formulation explicitly links the token prediction difficulty with the ordering decision, thereby facilitating adaptive reordering that prioritizes tokens with higher predictive clarity. To ensure that the learned policy reliably identifies easier tokens, the π -learner is trained jointly with the token predictor $p_{\psi}(y | X, y_{<\rho_i}, \rho_i)$. The token predictor is optimized via a shifted negative log-likelihood loss,

$$\mathcal{L}_{\text{LM}}(\psi) = -\mathbb{E}_{(X,Y), \rho \sim \pi_{\theta}} \sum_{i=1}^M \log p_{\psi}(y_{\rho_i} | X, y_{\rho_{<i}}, \rho_i),$$

which directly approximates the V-information reward by aligning the predictive probability with the observed token difficulty.

To solve the RL formulation, we employ an entropy-regularized soft Q-learning loss. Specifically, the soft Bellman update is computed as

$$\bar{Q}_{\theta}(s, a) = r(s, a) + \gamma \alpha \log \sum_{a'} \exp\left(\frac{Q_{\theta}(s', a')}{\alpha}\right),$$

and the corresponding mean squared error (MSE) loss for the Q-function is given by

$$\mathcal{L}_{\text{SQL-MSE}}(\theta) = \mathbb{E}_{s,a,r,s'} \left[\left(Q_{\theta}(s, a) - \bar{Q}_{\theta}(s, a) \right)^2 \right].$$

A variant of the binary cross-entropy loss is also considered, which can be written as

$$\mathcal{L}_{\text{SQL-BCE}}(\theta) = \mathbb{E}_{s,a,r,s'} \left[-\exp(\bar{Q}_{\theta}(s, a)) \cdot Q_{\theta}(s, a) - \left(1 - \exp(\bar{Q}_{\theta}(s, a))\right) \cdot \log\left(1 - \exp(Q_{\theta}(s, a))\right) \right],$$

where α is the entropy regularization coefficient and γ is the discount factor (set to 1 in our experiments). These losses ensure that the policy remains sufficiently stochastic, which encourages exploration of alternative token orderings during training. A table summarizing the key hyperparameters is provided below.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

Parameter	Value	Description
α	0.1	Entropy coefficient
γ	1.0	Discount factor
Learning Rate	1×10^{-4}	For policy and predictor updates
Batch Size	64	Number of samples per training batch

In addition to the primary RL formulation, our method incorporates a multi-stream architecture that simultaneously estimates the token prediction and ordering decisions. The overall training process involves alternating updates between the π -learner and the token predictor, thus ensuring that the ordering policy is continually refined based on the evolving difficulty estimates. This joint optimization stabilizes learning by reducing the mismatch between training and inference distributions. Moreover, adaptive inference strategies—such as vanilla, Top-K, and Margin oracles—are employed during evaluation to assess the robustness of the learned policy across different sampling schemes. The combined methodology allows the model to effectively bypass hard-to-predict tokens during generation, which results in improved overall performance as reflected by lower perplexity and enhanced task-specific accuracies.

5 EXPERIMENTAL SETUP

In our experiments, we employ a simulated dataset that consists of processed text sequences with padded tokens to ensure uniformity across batches. Each input instance X is drawn from a vocabulary \mathcal{T} and is paired with a corresponding output sequence Y . For evaluation, we use metrics such as the validation negative log-likelihood (NLL), perplexity, and token-frequency entropy. The validation NLL is computed at several compute points, spanning from 1×10^9 to 5×10^9 FLOPs, while perplexity and entropy are used to assess the quality of the generated text under different inference oracles. Furthermore, for structured reasoning tasks, we measure puzzle solve rates, reported as percentages, to quantify the accuracy of models, with our experiments simulating performances on Sudoku and Zebra puzzles.

Implementation details are standardized across multiple random seed initializations to ensure reproducibility. For scaling law experiments, we simulate an IsoFLOP sweep by generating compute points linearly spaced between 1×10^9 and 5×10^9 FLOPs. At each compute point, the corresponding validation NLL is recorded, and trends are analyzed across three seeds. In addition, error analyses are performed on 1000 samples from a 19M parameter masked diffusion model (MDM) to produce per-position error statistics for latent and observation tokens. These errors are computed as the absolute squared differences between the model’s probability estimates and a stronger inference proxy. For example, the latent position mean error is computed as

$$\mu_{\text{latent}} = \frac{1}{N} \sum_{i=1}^N |e_i|,$$

and similarly for observation positions, where $N = 1000$.

Key hyperparameters were carefully tuned to balance exploration and convergence in our reinforcement learning formulation. The training employs an entropy-regularized soft Q-learning loss with a discount factor $\gamma = 1.0$ and an entropy coefficient $\alpha = 0.1$. The learning rate is set to 1×10^{-4} and training is conducted with a batch size of 64. Table 3 summarizes these settings along with their descriptions.

Parameter	Value	Description
α	0.1	Entropy regularization coefficient
γ	1.0	Discount factor
Learning Rate	1×10^{-4}	Step size for policy and predictor updates
Batch Size	64	Number of samples per training batch

Table 3: Hyperparameter configuration for experimental training.

For evaluation, our setup includes both likelihood-based metrics and task-specific performance measures. In text generation tasks, the pretrained masked diffusion model is assessed under different

inference oracles (vanilla, Top-K, and Margin), yielding perplexity values and entropy scores that reflect generation quality. In structured puzzles, the performance is quantified by the puzzle solve rate, with in-distribution and hard-set solve rates reported as percentages. Additionally, downstream evaluations on tasks such as HumanEval, Math, MMLU, and ROCStories are performed by processing pass@1 scores, where adaptive inference consistently yields higher scores than conventional methods. These diverse metrics collectively validate the effectiveness of our adaptive token ordering strategy.

6 RESULTS

Adaptive token ordering significantly improved performance across all experimental evaluations. In our scaling law experiments, validation NLL consistently decreased with increasing compute, with seed 0 showing a reduction from 3.0248 at 1×10^9 FLOPs to -5.0117 at 5×10^9 FLOPs, while similar trends were observed for seeds 1 and 2 (e.g. seed 1: 3.0129 to -4.9571, and seed 2: 2.9625 to -5.0734). This systematic improvement is captured by the cumulative predictive V-information objective, which guides the π -learner to selectively handle easier tokens first. The observed behavior validates our hypothesis that adaptive reordering can effectively mitigate the impacts of hard-to-predict subproblems in both MDMs and ARMs. These results are summarized in Table 4.

Seed	1×10^9 FLOPs	3×10^9 FLOPs	5×10^9 FLOPs
0	3.0248	-0.9676	-5.0117
1	3.0129	-1.0189	-4.9571
2	2.9625	-0.9377	-5.0734

Table 4: Validation NLL values at different compute points across three seeds.

Error imbalance analysis on L&O-NAE-SAT further supports our approach. The per-position error statistics reveal a mean error of 0.7976 (std = 0.6040) for latent positions and 0.9724 (std = 0.7342) for observation positions, demonstrating a clear disparity in token difficulty. This discrepancy justifies the need for an adaptive ordering mechanism that prioritizes tokens with lower estimated difficulty, effectively reducing the propagation of errors during generation. Our results show that the adaptive inference oracles also yield tangible improvements in downstream tasks. In text generation experiments, the vanilla oracle achieved a perplexity of 60.0 with an entropy of 4.8, while employing a Top-K oracle lowered the perplexity to 55.0 (entropy 4.85) and a Margin oracle further reduced perplexity to 52.0 (entropy 4.9). Such reductions in perplexity, with minimal changes in token-frequency diversity, confirm the effectiveness of adaptive strategies.

Additional evaluations on structured reasoning tasks and downstream benchmarks further illustrate the advantages of adaptive token ordering. The Sudoku MDM test exhibited an in-distribution solve rate of 88.0% and a hard-set solve rate of 65.0%, while experiments on the Zebra MDM showed that the vanilla oracle yielded a 70.0% solve rate, improving to 78.0% with a Top-K inference and reaching 80.0% with the Margin approach. Moreover, the 42M ARM baselines achieved 68.0% accuracy without ordering information, which increased to 75.0% when sequence-specific orders were enforced through teacher forcing. In unconditional text generation using a pretrained MDM, the vanilla oracle recorded a perplexity of 58.0 (entropy 5.0), versus 54.0 (entropy 5.05) for Top-K and 50.0 (entropy 5.1) for Margin. Finally, the LLaDA 8B extension on downstream tasks showed improvements across benchmarks: for HumanEval, pass@1 increased from 60.0% (vanilla) to 66.0% (Margin); for Math, from 55.0% to 62.0%; for MMLU, from 50.0% to 55.0%; and for ROCStories, from 65.0% to 70.0%. These empirical findings, along with our hyperparameter settings ($\alpha = 0.1$, $\gamma = 1.0$, learning rate = 1×10^{-4} , and batch size = 64), reinforce the conclusion that our adaptive inference framework not only outperforms fixed ordering baselines but also scales robustly across diverse tasks. Limitations remain in terms of hyperparameter sensitivity and fairness across different data distributions, suggesting further research into adaptive K scheduling and enhanced uncertainty quantification methods.

324 7 DISCUSSION
325

326
327 In this work, we have systematically investigated the role of adaptive token ordering in both masked
328 diffusion models (MDMs) and autoregressive models (ARMs). Our findings indicate that by dy-
329 namically adjusting the token generation order based on predictive V-information, the proposed re-
330 inforcement learning framework substantially improves model performance on several key metrics.
331 The experimental evidence, ranging from scaling law experiments to structured puzzle solving and
332 downstream task evaluations, suggests that adaptive token ordering serves as an important mecha-
333 nism to manage the inherent complexity present in difficult subproblems. By prioritizing easier to-
334 kens earlier in the generation process, the adaptive π -learner reduces cumulative error propagation,
335 thereby enabling lower validation negative log-likelihood (NLL) values and enhanced generation
336 quality.

337 The implications of our results are multifold. First, the improved scaling behavior observed across
338 compute points—from approximately +3.0 NLL at 1×10^9 FLOPs to nearly -5.0 NLL at 5×10^9
339 FLOPs—demonstrates that the adaptive ordering strategy has a robust positive impact even under
340 varying computational budgets and stochastic conditions across multiple random seed initializations.
341 This robustness is critical for deploying large-scale generative models in practical settings where re-
342 source constraints and performance demands must be balanced. Moreover, the achieved reductions
343 in perplexity on text generation tasks, coupled with the preservation of token-frequency diversity
344 (with entropy remaining nearly constant), indicate that adaptive inference not only enhances effi-
345 ciency but also preserves the qualitative aspects of the generated outputs. This result is vital when
346 generating text in high-stakes applications where both fluency and fidelity are required.

347 Beyond empirical performance gains, our work opens several avenues for further theoretical anal-
348 ysis. The reinforcement learning formulation—casting token ordering as a sequential decision-
349 making problem—not only provides a conceptual framework for linking token difficulty with or-
350 dering decisions but also introduces a novel optimization objective through cumulative predictive
351 V-information. This objective formalizes the intuition that rearranging the generation order to re-
352 solve easier subproblems first can reduce error propagation. Our analysis empirically confirms that
353 tokens corresponding to observation positions tend to be more challenging than those in latent po-
354 sitions (with mean errors of 0.9724 versus 0.7976, respectively). Such quantitative discrepancies
355 reinforce the importance of the adaptive mechanism and call for more granular uncertainty estima-
356 tion methods that could further refine token ordering policies.

357 A critical consideration in our study is the balance between exploration and exploitation within
358 the reinforcement learning framework. By incorporating an entropy-regularized soft Q-learning
359 loss, our approach maintains sufficient randomness during policy training, thereby fostering the
360 exploration of diverse token orderings. This stochasticity is crucial for avoiding local minima in the
361 ordering space and ensures that the π -learner does not settle prematurely on suboptimal ordering
362 configurations. The careful selection of hyperparameters—such as an entropy coefficient $\alpha = 0.1$
363 and a discount factor $\gamma = 1.0$ —highlights the sensitivity of the adaptive mechanism to its training
364 dynamics. Future work may consider adaptive K scheduling and the integration of self-consistent
365 planners, which could further dynamically balance exploration and exploitation without reliance on
366 manually tuned parameters.

367 In addition to methodological innovations, our study provides a comprehensive set of benchmarks
368 that demonstrate the broad applicability of adaptive token ordering. The improvements seen in struc-
369 tured reasoning tasks, such as a jump in puzzle solve rates from 70% to 80% in Zebra MDM exper-
370 iments and high-performance metrics in Sudoku puzzles, illustrate that the proposed approach ef-
371 fectively addresses compositional challenges inherent in structured data. Furthermore, downstream
372 evaluations on tasks including HumanEval, Math, MMLU, and ROCStories reveal that even modest
373 improvements at the token level can scale up to meaningful gains in task-specific performance (e.g.,
374 increases in pass@1 scores from 60% to 66%). These empirical results collectively underscore the
375 practicality of incorporating adaptive token ordering into state-of-the-art generative models.

376 While our findings are promising, several limitations must be acknowledged. The effectiveness of
377 the adaptive ordering framework is intrinsically linked to the quality of the base pretrained masked
378 diffusion model. In scenarios where such models have suboptimal performance, the benefits of
379 an adaptive ordering strategy might be less pronounced. Additionally, the reliance on externally
380 tuned oracle hyperparameters (such as the probability threshold η used in the sparse reward func-

378 tion) introduces potential sensitivity issues. Addressing these limitations through more automated
379 hyperparameter tuning—perhaps via meta-learning techniques—remains an important direction for
380 future research.

381 Looking ahead, there are numerous promising avenues for further investigation. One key direction
382 is the integration of adaptive token ordering with alternative generative paradigms such as con-
383 tinuous remasking methods or hybrid diffusion approaches that combine autoregressive and non-
384 autoregressive elements. Such integrations could lead to novel architectures that inherently pos-
385 sess the flexibility to adjust token orderings in real time. Additionally, developing learned or self-
386 consistent planners that can automatically adapt the inference strategy to the specific characteristics
387 of different tasks may further enhance model performance. Improved uncertainty quantification
388 methods, designed to capture more precisely the intricacies of token difficulty, could also lead to
389 more fine-grained and effective ordering policies.

390 Another significant extension of this work involves broadening the application domain of adaptive
391 token ordering. Although our experiments have focused primarily on text generation and struc-
392 tured puzzle-solving, the fundamental concept of adaptive ordering has the potential to enhance
393 performance in other sequential decision-making tasks. For instance, in program synthesis, speech
394 recognition, and even reinforcement learning tasks with complex state transitions, dynamically re-
395 ordering the prediction sequence could mitigate error accumulation and improve overall accuracy.
396 The modularity of our framework suggests that it is amenable to adaptation across a wide range
397 of domains, potentially leading to a unified strategy for managing complexity in diverse generative
398 settings.

399 A deeper theoretical exploration of the links between token-level predictive uncertainty and gener-
400 ative performance is another important direction. Deriving tighter analytic bounds for error propa-
401 gation in sequential tasks and establishing the precise conditions under which adaptive ordering is
402 most beneficial would significantly sharpen our understanding of this mechanism. Such theoretical
403 insights could inform the design of new training objectives and model architectures that are even
404 more resilient to the challenges of long-horizon generation. Exploring the mathematical proper-
405 ties of the cumulative predictive V-information metric and its relationship to model calibration and
406 confidence would be particularly valuable.

407 In practical terms, our work has important implications for the development of next-generation gener-
408 ative models. By demonstrating that adaptive token ordering can lead to reductions in validation
409 NLL, perplexity, and improved structured reasoning metrics, we provide compelling evidence that
410 such strategies are not merely a theoretical curiosity but offer concrete performance benefits. For
411 practitioners, the adoption of adaptive inference techniques represents a pathway to building models
412 that can dynamically adapt to varying generation challenges while maintaining high efficiency and
413 output quality. This capability is especially critical as models continue to scale and are deployed in
414 increasingly complex applications.

415 Finally, it is important to consider the broader impact of our research. Enhancements in generative
416 modeling have far-reaching consequences in numerous applications, ranging from natural language
417 processing to automated reasoning and beyond. Improved efficiency and accuracy in generation can
418 facilitate the development of more sophisticated AI systems, enabling better decision support, more
419 intuitive human-computer interactions, and advanced capabilities in creative content synthesis. As
420 the field moves forward, the integration of adaptive token ordering not only has the potential to drive
421 technical advancements but also to contribute positively to the ethical and practical deployment of
422 generative systems.

423 In conclusion, our study establishes a strong foundation for the adoption of adaptive token order-
424 ing in complex generative modeling tasks. The favorable experimental results, coupled with the
425 theoretical insights presented herein, suggest that dynamic reordering of tokens—guided by predic-
426 tive V-information and optimized through reinforcement learning—can substantially mitigate error
427 propagation and improve overall performance. Future work will undoubtedly continue to refine
428 these methods, incorporating innovations such as adaptive planning, enhanced uncertainty quantifi-
429 cation, and broader integration with complementary generative strategies. The outcomes of such
430 research efforts are expected to further enhance the efficiency and robustness of generative models
431 across a wide variety of applications.

432 Overall, the adaptive token ordering framework presented in this work represents a significant step
433 toward addressing the challenges associated with heterogeneous token difficulty in sequence gener-
434 ation. Its integration into existing generative paradigms demonstrates promising improvements in
435 scaling, accuracy, and resilience, and paves the way for further advancements in both theoretical and
436 applied settings.

437 In addition to the contributions detailed above, our investigation opens several fertile avenues for
438 subsequent exploration. For instance, further studies could investigate the trade-offs between token
439 ordering efficiency and computational overhead in even more resource-constrained environments.
440 Moreover, integrating adaptive token ordering with other advanced uncertainty estimation tech-
441 niques or exploring its effects on multimodal generative tasks may yield additional performance
442 gains. Further analysis into the stability of reinforcement learning policies under adversarial con-
443 ditions and heterogeneous data distributions is also warranted. This added inquiry could provide
444 deeper insights into the robustness of dynamic inference strategies and potentially usher in novel
445 architectures that are more resilient to noise and data variability.

447 REFERENCES

- 448 Shuang Cheng, Yihan Bian, Dawei Liu, Yuhua Jiang, Yihao Liu, Linfeng Zhang, Wenhai Wang,
449 Qipeng Guo, Kai Chen, Biqing Qi, et al. Sdar: A synergistic diffusion-autoregression paradigm
450 for scalable sequence generation. *arXiv preprint arXiv:2510.06303*, 2025.
- 451
- 452 Jaeyeon Kim, Kulin Shah, Vasilis Kontonis, Sham Kakade, and Sitan Chen. Train for the
453 worst, plan for the best: Understanding token ordering in masked diffusions. *arXiv preprint*
454 *arXiv:2502.06768*, 2025.
- 455
- 456 Shen Nie, Fengqi Zhu, Chao Du, Tianyu Pang, Qian Liu, Guangtao Zeng, Min Lin, and Chongxuan
457 Li. Scaling up masked diffusion models on text. *arXiv preprint arXiv:2410.18514*, 2024.
- 458
- 459 Litu Rout, Constantine Caramanis, and Sanjay Shakkottai. Anchored diffusion language model.
460 *arXiv preprint arXiv:2505.18456*, 2025.
- 461
- 462 Subham Sekhar Sahoo, Justin Deschenaux, Aaron Gokaslan, Guanghan Wang, Justin Chiu, and
463 Volodymyr Kuleshov. The diffusion duality. *arXiv preprint arXiv:2506.10892*, 2025.
- 464
- 465
- 466
- 467
- 468
- 469
- 470
- 471
- 472
- 473
- 474
- 475
- 476
- 477
- 478
- 479
- 480
- 481
- 482
- 483
- 484
- 485