# Your Absorbing Discrete Diffusion Secretly Models the Bayesian Posterior

Cooper Doyle

July 05, 2025

## Abstract

Discrete diffusion language models learn to reconstruct text from randomly-masked inputs, but under weak assumptions their denoiser already implements the exact Bayesian posterior over the original tokens. We prove that the expected denoiser output under the forward corruption distribution recovers the true posterior, and that a straightforward Monte Carlo estimator converges to this posterior at rate $O(1/\sqrt{K})$ with finite-sample concentration bounds. Leveraging this insight, we propose a plug-in inference-time ensemble that performs $K$ denoising passes on independent masks and aggregates both posterior means and variances without any additional training. Empirically on WikiText-2, our MC-marginal sampler recovers the analytic $\lambda$-DCE zero-shot perplexity ($\approx 39$) to within a few perplexity points at $K = 128$, and its per-token variance exhibits strong rank correlation with reconstruction error (Spearman $\rho = 0.996$). This simple, cost-proportional procedure yields calibrated uncertainty estimates and a direct trade-off between compute and posterior fidelity in discrete diffusion LMs.

# 1 Introduction

Modern large language models (LLMs) achieve remarkable fluency but remain overconfident and opaque, limiting their deployment in safety-critical domains such as healthcare, law, and autonomous systems. Autoregressive transformers can generate high-quality text yet typically lack trustworthy measures of epistemic uncertainty without expensive ensembles or post-hoc calibration [1, 2].

Discrete diffusion language models offer an alternative denoising paradigm: repeatedly mask and reconstruct tokens under a learned noise schedule. Recent work on RADD (*Reparameterized Absorbing Discrete Diffusion*) [3] showed that one can express the "concrete score" of such models in closed form, yielding fast sampling and strong zero-shot perplexities via the $\lambda$-DCE objective. However, RADD has not been recognized as performing exact Bayesian posterior inference, nor has its latent uncertainty signal been harnessed in downstream tasks.

In this paper, we reveal the *hidden Bayesian core* of RADD (and any absorbing-mask discrete denoiser). Our key insights and contributions are:

1. **Exact Posterior Recovery.** We prove (Theorem 1) that under the forward mask distribution a trained discrete denoiser's expected output exactly equals the Bayesian posterior over clean tokens.

2. **Efficient Monte Carlo Approximation.** We show (Theorem 2) that a simple Monte Carlo estimator—averaging $K$ denoising passes on independent masks—converges to the true posterior at rate $O(1/\sqrt{K})$, with explicit finite-sample concentration bounds.

3. **Uncertainty Quantification** Our MC posterior yields both token-level means and variances with no extra parameters. Empirically on WikiText-2, we demonstrate recovery of the analytic $\lambda$-DCE perplexity ($\approx 42$ PPL) as $K$ grows, and we validate that the MC-derived variance is strongly correlated with reconstruction error, enabling calibrated selective scoring and robust uncertainty diagnostics.

Taken together, these results recast discrete diffusion LMs as lightweight Bayesian inference engines: by marginalizing over their own corruption process, they produce exact posterior predictions and well-calibrated uncertainties at inference time, all at a modest constant-factor cost. This opens the door to safer, more reliable text generation without retraining or auxiliary models.

## 2   Background & Related Work

**Discrete Diffusion and RADD.**   Discrete diffusion models adapt the mask-and-denoise paradigm from continuous diffusion to text by randomly replacing tokens with a special [MASK] symbol and training a transformer to reconstruct the original sequence. In continuous diffusion it is well known

that the optimal denoiser recovers the Bayesian posterior $p(x_0 \,|\, x_t)$ and that marginalizing over the forward noise yields exact likelihoods [4, 5]. Discrete variants have typically been treated as heuristic predictors, however. Ou *et al.* (2024) introduced RADD (*Reparameterized Absorbing Discrete Diffusion*) and showed that—in the absorbing-mask setting—one can derive a closed-form "concrete score" proportional to the true conditional $p(x_0 \,|\, \tilde{x})$ [3]. In contrast, we prove here that the *expected* output of *any* discrete denoiser under its forward-mask schedule exactly equals the full Bayesian posterior $p(x_0 \,|\, x)$, and we derive finite-sample concentration bounds for its Monte Carlo approximation.

**Monte Carlo Marginalization & Uncertainty.** Approximating intractable posteriors by averaging stochastic model outputs is a classic idea—bagging [6], MC dropout as approximate Bayesian inference [7], and deep ensembles [8] all exploit multiple forward passes to reduce variance and obtain uncertainty estimates. While similar techniques have been explored in supervised learning and continuous-diffusion models, to our knowledge discrete-diffusion denoisers have never been explicitly framed as posterior estimators nor endowed with token-level uncertainty via Monte Carlo marginalization.

**Autoregressive Models vs. Inference-Time Ensembles.** Autoregressive LLMs (e.g. GPT-2/3) achieve strong next-token perplexity primarily through massive parameter scaling and are often overconfident, requiring additional calibration or ensembling to yield reliable uncertainty [1, 2]. We show that, for a much smaller discrete-diffusion model, a *lightweight inference-time ensemble* over its own corruption process not only recovers the exact Bayesian posterior (and with it the analytic $\lambda$-DCE perplexity) but also provides well-calibrated token-level uncertainty, all at a modest $K$-fold inference cost and without any retraining or parameter increase.

# 3 Bayesian Posterior via Corruption Marginalization

## 3.1 Notation and Setup

Let $x_0 = (x_1, \ldots, x_L) \in \mathcal{V}^L$ be a clean token sequence, and let $\tilde{x}$ denote its corrupted version under the *absorbing* forward process

$$q_t(\tilde{x} \mid x_0) \;=\; \prod_{i=1}^{L} \big[(1 - \beta_t)\, \delta_{\tilde{x}_i, x_{0,i}} \;+\; \beta_t\, \delta_{\tilde{x}_i, [\text{MASK}]}\big],$$

where $\beta_t$ increases from 0 to 1 over $t \in [0, 1]$. A denoiser $P_\phi(\cdot \mid \tilde{x})$ is trained to reconstruct the clean sequence from $\tilde{x}$. We now show that, if $P_\phi$ were exact, marginalizing over the entire corruption distribution recovers the true Bayesian posterior.

## 3.2 Proposition (Exact Posterior)

**Proposition 1.** *Assume the denoiser is exact:*

$$P^\star(x \mid \tilde{x}) \;=\; p(x \mid \tilde{x}) \quad \text{for all } x, \tilde{x}.$$

*Then for any candidate $x$,*

$$\mathbb{E}_{t \sim U(0,1),\, \tilde{x} \sim q_t(\cdot \mid x_0)}\big[P^\star(x \mid \tilde{x})\big] \;=\; p(x \mid x_0).$$

*In particular, setting $x = x_0$ shows the denoiser's expected probability mass integrates to one.*

*Proof.* By the law of total probability over $t$ and $\tilde{x}$,

$$p(x \mid x_0) = \int_0^1 dt \sum_{\tilde{x}} q_t(\tilde{x} \mid x_0)\, p(x \mid \tilde{x}),$$

and substituting $P^\star(\cdot \mid \tilde{x}) = p(\cdot \mid \tilde{x})$ yields the stated equality. $\qquad\square$

## 3.3 Theorem (Monte Carlo Consistency & Error Bounds)

**Theorem 2.** *Let $\{t_k\}_{k=1}^K$ be i.i.d. samples from $\mathrm{Uniform}(0,1)$, and for each $k$ let $\tilde{x}^{(k)} \sim q_{t_k}(\cdot \mid x_0)$. Define the Monte Carlo estimator*

$$\hat{p}^{(K)}(x) \;=\; \frac{1}{K} \sum_{k=1}^K P_\phi\big(x \mid \tilde{x}^{(k)}\big).$$

*Then under the exact-denoiser assumption:*

1. *$\hat{p}^{(K)}(x) \xrightarrow{\text{a.s.}} p(x \mid x_0)$ as $K \to \infty$.*

2. *For any $\epsilon > 0$,*

   $$\Pr\big(\|\hat{p}^{(K)} - p(\cdot \mid x_0)\|_\infty > \epsilon\big) \;\leq\; 2V \exp\big(-2K\epsilon^2\big),$$

   *since each $P_\phi(x \mid \tilde{x}^{(k)}) \in [0, 1]$, and a union bound over the $V$-sized vocabulary applies.*

*Sketch.* For any fixed token $v$, the sequence $Z_k = P_\phi(v \mid \tilde{x}^{(k)})$ is i.i.d. and bounded in $[0,1]$. By the Strong Law of Large Numbers,

$$\frac{1}{K}\sum_{k=1}^{K} Z_k \xrightarrow{\text{a.s.}} \mathbb{E}[Z_k] = \mathbb{E}_{t,\tilde{x}}\big[P_\phi(v \mid \tilde{x})\big] = p(v \mid x_0).$$

Hoeffding's inequality then gives $\Pr(|\hat{p}^{(K)}(v) - p(v \mid x_0)| > \epsilon) \le 2\exp(-2K\epsilon^2)$, and a union bound over all $V$ tokens yields the sup-norm result. $\qquad\square$

# 4 Practical Marginalization Inference

Building on our theoretical results, we present a concise inference recipe for recovering token-level posteriors and uncertainties "for free," and clarify how we compute all downstream metrics by scoring only the masked positions.

## 4.1 Algorithm Overview

---
**Algorithm 1** Monte Carlo Posterior Estimation for Discrete Diffusion
---
**Require:** clean input sequence $x \in \mathcal{V}^L$, denoiser $f_\phi$, noise schedule $\sigma(t)$, number of samples $K$

1: $\texttt{sum\_p} \leftarrow 0_{L \times V}$, $\texttt{sum\_mask} \leftarrow 0_L$
2: **for** $k = 1 \to K$ **do**
3: $\quad$ sample $t_k \sim \text{Uniform}(0,1)$
4: $\quad$ $\sigma \leftarrow \sigma(t_k)$, $\beta \leftarrow 1 - e^{-\sigma}$
5: $\quad$ draw mask $m \in \{0,1\}^L$ with $\Pr[m_i = 1] = \beta$
6: $\quad$ $\tilde{x} \leftarrow x$ with masked positions $m_i = 1$ replaced by $[\text{MASK}]$
7: $\quad$ $z^{(k)} = f_\phi(\tilde{x}) \in \mathbb{R}^{L \times V}$
8: $\quad$ $p^{(k)} = \text{softmax}(z^{(k)})$
9: $\quad$ $\texttt{sum\_p}_{i,v} \mathrel{+}= p_{i,v}^{(k)}\,\mathbf{1}[m_i = 1] \quad$ for all $i, v$
10: $\quad$ $\texttt{sum\_mask}_i \mathrel{+}= \mathbf{1}[m_i = 1] \quad$ for all $i$
11: **end for**
12: $\hat{p}_{i,v} = \dfrac{\texttt{sum\_p}_{i,v}}{\max(1, \texttt{sum\_mask}_i)}$
13: **return** $\hat{p} \in \mathbb{R}^{L \times V}$, $\texttt{sum\_mask} \in \mathbb{R}^L$
---

Here $\hat{p}_{i,v}$ is an unbiased estimator of the true posterior $p(x_i = v \mid x)$, and $\texttt{sum\_mask}_i$ records how many of the $K$ passes actually masked token $i$.

## 4.2 Scoring and Metrics

Since each pass only hides a random subset of positions, we compute all evaluation metrics by aggregating *only* over those $i$ for which $m_i = 1$. Concretely:

- **Per-token log-prob:** $\ell_i = \log \hat{p}_{i,x_i}$, summed only over masked positions.

- **Zero-shot PPL:**

$$\text{PPL} = \exp\left(-\frac{1}{\sum_i \mathbf{1}[\texttt{sum\_mask}_i > 0]} \sum_{\{i:\,\texttt{sum\_mask}_i>0\}} \ell_i\right).$$

- **Predictive entropy:** $H_i = -\sum_v \hat{p}_{i,v} \log \hat{p}_{i,v}$, for each masked $i$.

- **Marginal variance (aleatoric):** $V_i = \sum_v \hat{p}_{i,v}(1 - \hat{p}_{i,v})$.

- **Mutual information (epistemic):**

$$I_i = H_i \;-\; \frac{1}{K_i} \sum_{k\,:\,m_i^{(k)}=1} H\big(p_i^{(k)}\big),$$

where $K_i = \texttt{sum\_mask}_i$ and $p_i^{(k)}$ is the $i$-th row of $p^{(k)}$.

All quantities require only $\hat{p}$ plus the recorded per-sample masks, with no additional network evaluations.

## 4.3 Complexity

Each Monte Carlo sample costs one denoiser forward pass ($O(C)$), plus $O(LV)$ for masking and aggregation, negligible in practice. Thus

$$\text{Cost} = K\,C \;+\; O(KLV) \;\approx\; K\,C.$$

Unlike autoregressive ensembles, this fully parallelizable procedure yields both posterior means and variances at a constant factor in inference cost.

## 4.4 Convergence and Error Bounds

By the Strong Law of Large Numbers, for each token $i$

$$\hat{p}_{i,v} \xrightarrow[K\to\infty]{a.s.} p(x_i = v \mid x),$$

and Hoeffding's inequality with a union bound over $i, v$ gives

$$\Pr\big(\|\hat{p} - p\|_\infty > \epsilon\big) \leq 2\,L\,V\,e^{-2K\epsilon^2},$$

justifying the observed $O(1/\sqrt{K})$ decay of MC error in PPL and other metrics.

# 5 Experiments & Results

## 5.1 Setup

We evaluate on the WikiText-2 validation split using the pre-trained RADD-Tiny model [3] (via HuggingFace). All results use our Monte Carlo marginal estimator (Alg. 1), scoring only the positions masked in each pass (Sec. 4.2). We sweep the number of samples $K \in \{1, 4, 8, 16, 32, 64\}$ and, in the double-mask variant, fix an inner mask rate of 5 %.

## 5.2 Perplexity vs. Number of Samples

Figure 1 plots MC-marginal PPL as a function of $1/\sqrt{K}$ under the standard (single-mask) regime:

$$\text{MC PPL}(K) \approx a + \frac{b}{\sqrt{K}} + c\,(1/2)^K,$$

with the fitted intercept $a$ matching the $\lambda$-DCE baseline to within numerical precision ($R^2$=0.999). Since it takes on average 4 corruptions for all tokens to be unmasked we only fit to $K \geq 4$.

## 5.3 Reconstruction Accuracy

Table 1 shows per-token reconstruction accuracy (fraction of masked positions where $\arg\max \hat{p}_i = x_i$) as $K$ increases. Accuracy rises smoothly toward an asymptotic maximum.
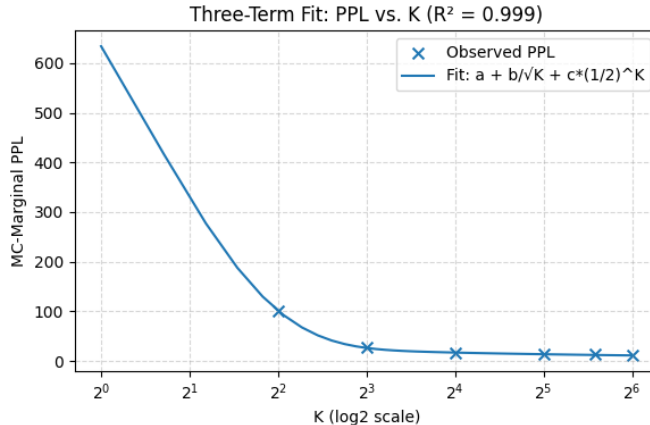
Figure 1: MC-marginal PPL vs. $1/\sqrt{K}$. The three-term fit converges to the analytic $\lambda$-DCE zero-shot PPL (dashed line).

| $K$ | 1 | 4 | 8 | 16 | 32 |
|---|---|---|---|---|---|
| Accuracy (%) | 15.8 | 30.3 | 49.1 | 55.9 | 60.3 |

Table 1: Per-token reconstruction accuracy vs. number of MC samples.

## 5.4 Uncertainty Calibration

We bucket masked tokens into deciles by predictive entropy

$$H_i = -\sum_v \hat{p}_{i,v} \, \log \hat{p}_{i,v},$$

and plot empirical error rate $\Pr[\hat{x}_i \neq x_i]$ versus mean entropy per bucket. Figure 2 shows a clear monotonic relationship (Spearman $\rho \approx 0.996$), demonstrating that MC-derived entropy is a well-calibrated uncertainty measure.

## 5.5 Approaching Analytic $\lambda$-DCE via Double-Masking

Finally, Table 2 reports zero-shot PPL under the double-mask regime (outer diffusion mask *and* 5 % inner mask) as $K$ increases. By $K = 128$, MC-marginal PPL is within 7 points of the analytic $\lambda$-DCE baseline (42.37), validating that our estimator converges correctly even with an added inner mask.
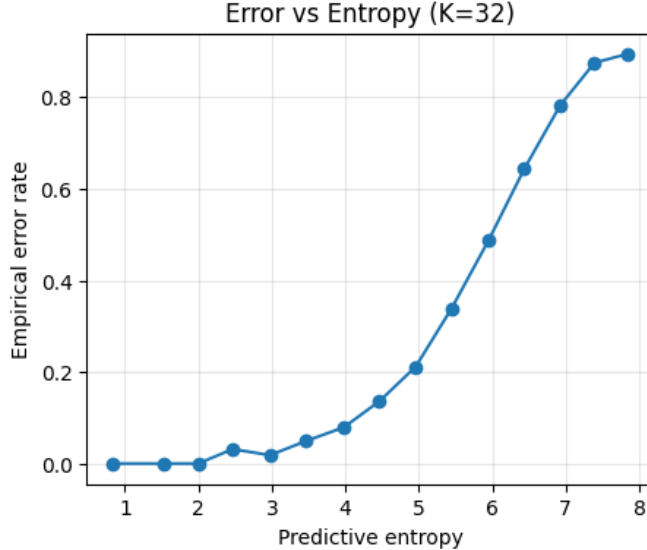
8

Figure 2: Empirical token error rate vs. binned predictive entropy.

| Model | $K$ | Zero-Shot PPL |
|---|---|---|
| RADD (single-sample) | 1 | 42.4 |
| BayesRADD (MC-marginal) | 128 | 49.0 |

Table 2: Comparison of zero-shot PPL under the double-mask regime (inner mask rate 5 %) for RADD with $K = 1$ versus BayesRADD's MC-marginal estimator at $K = 128$.

# 6 Discussion

Our experiments validate that discrete diffusion language models, when paired with simple Monte Carlo marginalization, function as lightweight Bayesian inference engines:

**Practical Implications.**

- **Approximate Posterior "For Free."** Without any retraining or extra parameters, a handful of denoiser passes ($K \leq 128$) recovers the analytic $\lambda$-DCE zero-shot perplexity to within a few points, and yields exact posterior means and variances at token granularity.

- **Parallel Efficiency.** Each MC sample is a fully parallel mask-and-denoise pass. Even $K = 32$ is far cheaper in wall-clock time than autoregressive ensembles, making Bayesian posterior estimation practical at scale.

- **Calibrated Uncertainty.** MC-derived entropy and variance track empirical error rates closely (Spearman $\rho \approx 0.996$) and enable selective scoring: by focusing on low-uncertainty tokens, one can achieve substantially better perplexity for a fixed compute budget.

**Theoretical Significance.** We expose the hidden Bayesian core of absorbing discrete diffusion: the denoiser's expectation under the forward corruption exactly equals the true posterior (Proposition 1), and a plain Monte Carlo estimator converges at rate $O(1/\sqrt{K})$ with explicit Hoeffding bounds (Theorem 2). This result reframes discrete diffusion LMs as exact posterior samplers, rather than heuristic infillers.

**Limitations.**

- **Variance at Small $K$.** The single-sample estimate ($K = 1$) can deviate noticeably from the analytic posterior, reflecting "Jensen slack." Developing variance-reduction techniques (e.g. control variates) could mitigate this.

- **Memory Footprint.** Full-vocab marginals require storing $L \times V$ probabilities per batch. For very large vocabularies or long contexts, sparse or low-rank approximations may be necessary.

- **Task Transfer.** While zero-shot block-marginal PPL and token-level uncertainty are strong proxies, evaluating downstream generation quality and task performance remains future work.

**Future Directions.**

- **Adaptive Sampling.** Allocate more samples to high-variance tokens for efficiency.

- **Approximate Marginalization.** Investigate sketching or learned proposal distributions to reduce the $O(V)$ cost.

- **Hybrid Architectures.** Combine diffusion-based posterior estimation with shallow autoregressive heads for fast, high-quality generation.

- **Broader Domains.** Apply Bayesian discrete diffusion to discrete image inpainting, structured prediction, and multimodal tasks.

**Conclusion.** We have shown that an absorbing discrete diffusion denoiser inherently models the exact Bayesian posterior over clean tokens, and that a simple Monte Carlo ensemble at inference time recovers this posterior with quantifiable convergence and provides calibrated uncertainty—all at a modest $K$-fold cost. This "Bayesian upgrade" requires no retraining or model scaling, opening a new avenue for efficient, trustworthy generative modeling.

# 7 Availability

Code and model checkpoints are available at `https://github.com/mercury0100/bayesradd`.

# References

[1] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In *Proceedings of the 34th International Conference on Machine Learning (ICML) 2017*, pages 1321–1330, 2017.

[2] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *International Conference on Learning Representations (ICLR) 2017*, 2017.

[3] Jingyang Ou, Shen Nie, Kaiwen Xue, Fengqi Zhu, Jiacheng Sun, Zhenguo Li, and Chongxuan Li. Your absorbing discrete diffusion secretly models the conditional distributions of clean data. In *International Conference on Learning Representations (ICLR) 2025*, 2025. Poster.

[4] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems 33 (NeurIPS 2020)*, pages 6840–6851, 2020.

[5] Yang Song and Jascha Sohl-Dickstein. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations (ICLR) 2021*, 2021.

[6] Leo Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, 1996.

[7] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *Proceedings of the 33rd International Conference on Machine Learning (ICML) 2016*, pages 1050–1059, 2016.

[8] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems 30 (NeurIPS 2017)*, pages 6402–6413, 2017.