
Any-Order Flexible Length Masked Diffusion

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Masked diffusion models (MDMs) have recently emerged as a promising alterna-
2 tive to autoregressive models over discrete domains. MDMs generate sequences in
3 an any-order, parallel fashion, enabling fast inference and strong performance on
4 non-causal tasks. However, a crucial limitation is that they do not support token
5 insertions and are thus limited to *fixed-length* generations. To this end, we introduce
6 **Flexible Masked Diffusion Models (FlexMDMs)**, a discrete diffusion paradigm
7 that simultaneously can model sequences of flexible length while provably retain-
8 ing MDMs’ flexibility of any-order inference. Grounded in an extension of
9 the stochastic interpolant framework, FlexMDMs generate sequences by inserting
10 mask tokens and unmasking them. Empirically, we show that FlexMDMs match
11 MDMs in perplexity while modeling length statistics with much higher fidelity.
12 On a synthetic maze planning task, they achieve $\approx 60\%$ higher success rate than
13 MDM baselines. Finally, we show pretrained MDMs can easily be *retrofitted* into
14 FlexMDMs: on 16 H100s, it takes only three days to fine-tune LLaDA-8B into a
15 FlexMDM, achieving superior performance on math (GSM8K, 58% \rightarrow 67%) and
16 code infilling performance (52% \rightarrow 65%).

17 1 Introduction

18 While diffusion models [Ho et al., 2020, Song et al., 2020, Sohl-Dickstein et al., 2015] are now the
19 leading paradigm for generative modeling in continuous domains, recent work has begun to expand
20 their scope to discrete spaces. The prevailing approach, Masked Diffusion Models (MDMs) [Shi
21 et al., 2024, Sahoo et al., 2024, Gat et al., 2024], generates sentences in a non-left-to-right, any-order
22 fashion. Compared to autoregressive models, this any-order generation ability yields substantially
23 faster inference and strong downstream performance on non-casual tasks such as planning [Ye et al.,
24 2024], code Nie et al. [2025], Ye et al. [2025], and reasoning [Nie et al., 2024].

25 Despite these successes, current MDMs cannot (1) model distributions supported on sequences of
26 *variable length* and (2) insert new tokens during generation (Figure 1, left). Both capabilities are
27 natural desiderata for generative models over discrete domains. We therefore ask: *Can we model*
28 *variable-length data while preserving MDMs’ any-order generation power?*

29 We answer in the affirmative by proposing the **Flexible Masked Diffusion Model (FlexMDM)**.
30 FlexMDMs start from an empty string and gradually [insert mask tokens](#) and then [unmask](#) them
31 (Figure 1, right). Beyond learning the usual *unmasking* posterior—the distribution of a clean token at
32 masked positions—we introduce an *insertion* expectation: the expected number of tokens to insert
33 conditioned on the current sequence. Crucially, we show that FlexMDM is [theoretically grounded](#)
34 (i.e., under perfect training, it samples from the true data distribution) and [retains the any-order](#)
35 [sampling property](#) of MDMs, thereby directly addressing the question above.

36 Empirically, we demonstrate that FlexMDM offers significant new upgrades to the MDM paradigm,

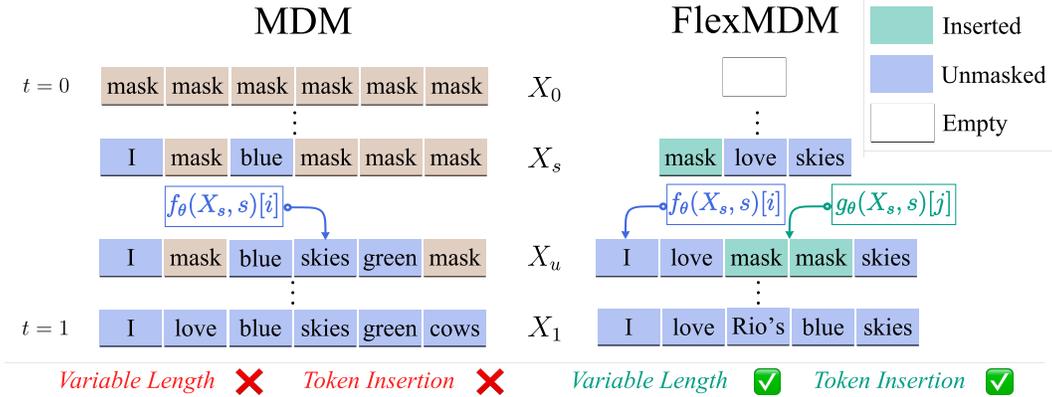


Figure 1: Flexible Masked Diffusion Model (FlexMDM) addresses MDMs’ inability to handle variable-length sequences and token insertion while preserving any-order generation power. At each step, FlexMDM performs **insertion** and **unmasking** by predicting **the expected number of mask tokens to insert** (g_θ) and the **posterior over clean tokens** (f_θ), respectively.

- 37 • A FlexMDM pretrained on OpenWebText is able to **model the length distribution with substantially**
- 38 **higher fidelity** while matching the perplexity of an MDM counterpart.
- 39 • On **planning tasks**, FlexMDM achieves markedly better results, beating the success rate of MDMs
- 40 by nearly 60% on a natural synthetic baseline.
- 41 • **MDMs can be retrofitted into FlexMDMs at 8B+ scale:** We fine-tune LLaDA-8B [Nie et al., 2025],
- 42 an open-source MDM, into a FlexMDM using only 16 H100s for three days. The model transfers
- 43 cleanly from its MDM initialization and, with its newly acquired variable-length capability, attains
- 44 notably better performance on GSM8K (58%→67%) and Code infilling (52%→65-%).

45 Theoretically, our construction relies on the machinery of continuous-time Markov chains (CTMCs)

46 and in particular introduces the new notion of a *joint interpolant*, a novel extension of *stochastic*

47 *interpolants* [Albergo and Vanden-Eijnden, 2022, Albergo et al., 2023a, Lipman et al., 2022]. Recent

48 work [Zheng et al., 2024, Ou et al., 2024] established an equivalence between MDMs and any-order

49 language models—obviating the need for CTMCs. In contrast, we prove that FlexMDMs also possess

50 the flexibility of any-order generation, yet the continuous-time perspective is absolutely essential for

51 them to accurately model the length distribution. Accordingly, we re-derive the connections between

52 MDMs and stochastic interpolants and use them to ground the design of FlexMDMs.

53 **Roadmap.** We begin in Section 2 with a broadly accessible review of CTMCs and the connection

54 between MDMs and discrete flow matching. Building on this, Section 3 derives the FlexMDM training

55 objective, and Section 4 introduces our inference procedures. Section 5 presents our experimental

56 results.

57 2 Preliminaries: Continuous-Time Markov Chains and Masked Diffusions

58 In what follows, we provide an overview of continuous-time Markov chains (CTMCs), their role in

59 defining discrete diffusion models, and link them to the MDM framework. As we mentioned in the

60 introduction, this theme is essential to defining FlexMDMs in Section 3.

61 **Transport with continuous-time Markov Chains.** Given a target distribution p_1 over sequences with

62 a finite vocabulary set (e.g., text), our aim is to learn to transport samples from a reference distribution

63 p_0 through a continuum of distributions $\{p_t\}_{t \in [0,1]}$ such that $p_{t=1} = p_1$. This type of transport can be

64 realized by a continuous-time Markov chain, which is a stochastic process $\{X_t\}_{t \in [0,1]}$ with $X_0 \sim p_0$

65 governed by a time-dependent transition rate matrix $\{R_t(\cdot, \cdot)\}_{t \in [0,1]}$ satisfying

$$R_t(x, x) = - \sum_{y \neq x} R_t(x, y), \quad R_t(x, y) \geq 0, \quad x \neq y. \quad (1)$$

66 Intuitively, the rate matrix determines the infinitesimal likelihood that X_t transitions to any other

67 state y via

$$\mathbb{P}(X_{t+h} = y | X_t = x) = \mathbf{1}_{\{x=y\}} + hR_t(x, y) + o(h), \quad (2)$$

68 where we denote the conditional probability measure $\mathbb{P}(\cdot | X_t = x)$ of a new state given the present

69 one. Here $o(h)$ is a remainder term that vanishes faster than h as $h \rightarrow 0$. In generative modeling for

70 these discrete distributions, our aims are to **(a)** specify a path of marginal distributions $\{p_t\}_{t \in [0,1]}$
 71 connecting p_0 to p_1 and **(b)** learn the associated R_t such that these marginals collectively satisfy the
 72 Kolmogorov forward equation:

$$\partial_t p_t(x) = \sum_y p_t(y) R_t(y, x) \quad p_{t=0} = p_0. \quad (3)$$

73 This ensures that at time $t = 1$, the evolution specified by (2) results in a sample from the target
 74 distribution p_1 . The rate matrices defined in this paper are sparse; therefore, we assume that the
 75 unspecified entries are 0 and the diagonal entries are defined through Equation (1).

76 2.1 Masked Diffusion Models

77 We briefly review MDMs [Sahoo et al., 2024, Shi et al., 2024] and discrete flow matching with
 78 the masked construction [Gat et al., 2024], through the lens of stochastic interpolants. The target
 79 distribution p_1 assigns probability to length L sequences. The base distribution p_0 employed by these
 80 models is the point mass distribution at the fully masked length- L sequence $(\mathbf{m}, \dots, \mathbf{m})$, where \mathbf{m} is
 81 an auxiliary mask token.

82 To define the intermediate $\{p_t\}_{t \in [0,1]}$ that bridges the base and the target, we make use of a [stochastic](#)
 83 [interpolant](#) $\{x_t\}_{t \in [0,1]}$, a collection of random variables whose marginal distribution defines the
 84 continuum $\{p_t\}_{t \in [0,1]}$, i.e., $x_t \sim p_t$. Although the previous notion of stochastic interpolant [Albergo
 85 et al., 2023a] is defined in a continuous space, it naturally extends to a discrete space, and we defer a
 86 formal exposition to Appendix C.

87 **Design of distribution path.** The stochastic interpolant relies on a smooth and monotone unmasking
 88 schedule $\alpha_t: [0, 1] \rightarrow [0, 1]$ with boundary condition $(\alpha_0, \alpha_1) = (0, 1)$ and time derivative denoted
 89 by $\dot{\alpha}_t$. To draw x_t , we first sample a clean sequence $x_1 \sim p_1$; then, independently for every coordinate
 90 i , we draw an unmasking time T^i from density $\dot{\alpha}_t dt$ and set

$$x_t^i = \begin{cases} \mathbf{m} & t < T^i \\ x_1^i & t \geq T^i \end{cases}.$$

91 Hence, each clean token stays masked with probability $1 - \alpha_t$ in a coordinate-independent fashion,
 92 defining $p_t(\cdot | x_1)$. We then write p_t by marginalizing over $x_1 \sim p_1$.

93 **MDM training.** We now derive the MDM rate matrix that induces a CTMC whose marginals
 94 coincide with $\{p_t\}_{t \in [0,1]}$ and how it is learned in practice. The central object is the [unmasking](#)
 95 [posterior](#): the posterior on the clean token x_1^i for masked index i given $x_t = x$ and time step t , i.e.,
 96 $\mathbb{P}(x_1^i = v | x_t = x)$. We model this posterior with a neural network $f_\theta(x, t) \in (\Delta(\Sigma))^n$, where $\Delta(\Sigma)$
 97 denotes a simplex of probability distributions over the vocabulary Σ .

98 For every position where $x^i = \mathbf{m}$, the network aims to predict $f_\theta(x, t)[i, v] \approx \mathbb{P}(x_1^i = v | x_t = x)$,
 99 and is trained by minimizing the following variational loss:

$$\mathcal{L}_\theta = - \int_0^1 \mathbb{E} \left[\frac{\dot{\alpha}_t}{1 - \alpha_t} \sum_{i: x_t^i = \mathbf{m}} \log f_\theta(x_t, t)[i, x_1^i] \right] dt. \quad (4)$$

100 Here, \mathbb{E} denotes the expectation over $x_1 \sim p_1$ and $x_t \sim p_t(\cdot | x_1)$. The minimizer of this loss is the
 101 ground-truth unmasking posterior, which fully determines the MDM's rate matrix below. Precisely,
 102 for $t \in [0, 1]$, the rate matrix at time t is given by: for a partially masked sequence $x \in (\Sigma \cup \{\mathbf{m}\})^L$,

$$R_t(x, x[x^i \leftarrow v]) = \frac{\dot{\alpha}_t}{1 - \alpha_t} \underbrace{\mathbb{P}(x_1^i = v | x_t = x)}_{\text{unmasking posterior}}, \quad v \in \Sigma, x^i = \mathbf{m}, \quad (5)$$

103 where $x[x^i \leftarrow v]$ denotes the sequence obtained from x by replacing its i -th token with v . Therefore,
 104 once f_θ has learned the unmasking posterior, one can simulate the CTMC using the rate matrix in
 105 (5). The variational loss (4) quantifies the sampling guarantee of this *estimated* CTMC. Let p_1^θ be the
 106 terminal distribution of the estimated CTMC. Then, the loss function bounds the KL-divergence:

$$\mathcal{D}_{\text{KL}}(p_1 || p_1^\theta) \leq \mathcal{L}_\theta - \mathcal{L}_*,$$

107 where \mathcal{L}_* is the global minimum of \mathcal{L} . When the loss is in its infimum, the KL divergence vanishes,
 108 resulting in the ground truth distribution.

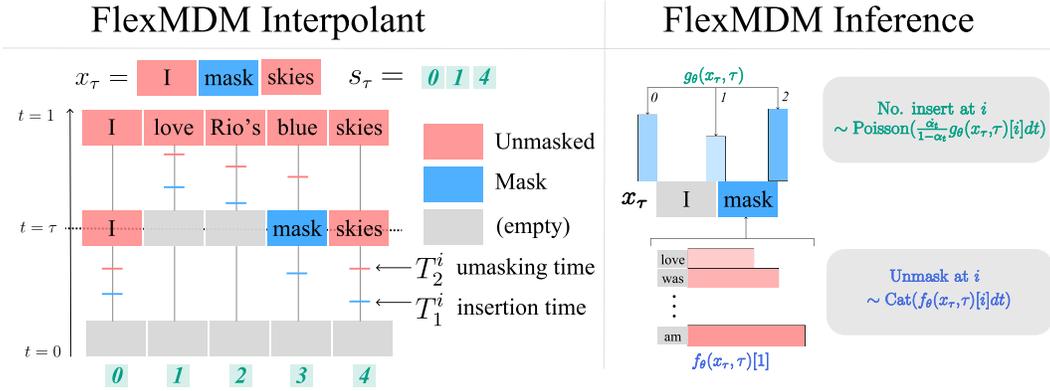


Figure 2: Left (FlexMDM interpolant). To draw a sample x_t , one can equivalently draw a sample $x_1 \sim p_1$, and for each token **unmask**, **mask**, or **remove** it according to the unmasking and insertion times (T_1^i, T_2^i). An auxiliary interpolant s_t gives closed-form expressions for the FlexMDM rate matrices. **Right (FlexMDM Inference).** We train the model to predict the [unmasking posterior](#) and the [insertion expectation](#) and demonstrate approximate inference strategy.

109 3 Variable Length Masked Diffusions: Training

110 In this section, we introduce **Flexible Length Masked Diffusion Model** (FlexMDM): a discrete
 111 diffusion that models a distribution p_1 assigning probabilities to sequences of different lengths.
 112 Following the MDM’s recipe, we aim to introduce a stochastic interpolant x_t whose marginal
 113 distribution defines the path $\{p_t\}_{t \in [0,1]}$ and learn the corresponding CTMC. Everything hinges on
 114 selecting an interpolant that is **(a)** easy to sample at $t = 0$ and **(b)** equipped with a closed-form rate
 115 matrix amenable to neural network training.

116 **Challenge.** Reusing the MDM interpolant is *inadequate*: at $t = 0$, the base distribution p_0 would
 117 consist of fully-masked sentences of *variable lengths*, which is impossible to sample since we do
 118 not know the length statistics of p_1 in advance. On the other hand, one can consider an interpolant
 119 constructed by masking and removing tokens from a clean sequence. However, this *complicates* the
 120 rate matrix characterization—token indices shift as insertion occur. To bridge this gap, we introduce
 121 the **joint interpolant**, an extension of the stochastic interpolant that augments the process with an
 122 auxiliary variable explicitly tracking token positions. This enlarged state space allows us to construct
 123 a broader class of rate matrices while preserving an easy-to-sample base distribution.

124 **Design of distribution path.** We now introduce our FlexMDM’s joint interpolant that allows us
 125 to model the variable length p_1 . This construction relies on *two* smooth, monotone schedules – an
 126 insertion schedule $\alpha: [0, 1] \rightarrow [0, 1]$ and an unmasking schedule $\beta: [0, 1] \rightarrow [0, 1]$, with the boundary
 127 conditions $(\alpha_0, \alpha_1) = (\beta_0, \beta_1) = (0, 1)$ and time derivatives denoted by $\dot{\alpha}_t, \dot{\beta}_t$.

128 To draw x_t , we first sample a clean sentence $x_1 \sim p_1$. Independently for each coordinate i , we
 129 draw an **insertion time** T_1^i and an **unmasking time** T_2^i with $T_1^i < T_2^i$ according to the density below.
 130 Accordingly, we either delete, mask, or unmask x_1^i to obtain x_t^i :

$$131 \quad T_1^i \sim \dot{\alpha}_t dt, \quad T_2^i \sim \mathbf{1}_{\{t \geq T_1^i\}} \frac{\dot{\beta}_t}{1 - \beta_{T_1^i}} dt, \quad x_t^i = \begin{cases} (\text{empty}), & 0 < t < T_1^i \\ \mathbf{m}, & T_1^i \leq t < T_2^i \\ x_1^i, & T_2^i \leq t \leq 1 \end{cases} \quad (6)$$

132 Here, $\mathbf{1}$ denotes the indicator function. We obtain x_t by concatenating the symbols x_t^i , and dropping
 133 $x_t^i = (\text{empty})$. Consequently, the length of x_t is equal to or less than that of x_1 ¹ (see Figure 2,
 134 left). As we mentioned above, we augment x_t with an index-tracking variable s_t , forming the joint
 interpolant (x_t, s_t) . Let $\text{len}(x_t)$ denote the length of x_t ; then

$$s_t := \{i \in \{1, \dots, \text{len}(x_1)\} \mid T_1^i \leq t\},$$

135 i.e., the set of indices whose clean tokens have *not* been deleted. Equivalently, the positions in x_1
 136 referenced by x_t ’s each index. By regarding s_t as a list and ordering its elements in ascending order,
 137 we also have $x_t = (x_1^{s_t[0]}, \dots, x_1^{s_t[\text{len}(s_t)-1]})$. We revisit s_t shortly to show how it enables an explicit

¹Writing x_t^i to mean the symbol derived from source position i is this a mild abuse of notation since the
 superscript i refers to a position in x_1 rather than a valid index of the (shorter) sequence x_t .

138 rate matrix. Since (x_t, s_t) is governed by the sampled unmasking and insertion times, we write
 139 $(x_t, s_t) \sim p_t(\cdot | x_1)$. Marginalizing $p_t(\cdot | x_1)$ over $x_1 \sim p_1$ yields p_t . Since the boundary condition
 140 sets $\alpha_0 = \beta_0 = 0$, all tokens are deleted at $t = 0$; thus, p_0 is the point mass on the empty string.

141 **FlexMDM training.** We now explain how we train our FlexMDM to learn the desired rate matrix.
 142 We first discuss what the CTMC looks like at a high level: recall from (6) that when t increases,
 143 tokens are progressively inserted and unmasked. Indeed, one can show that a CTMC that generates
 144 the interpolant can be characterized by two quantities that govern the rate of insertion and unmasking
 145 (see Figure 2, right):

- 146 • **Unmasking posterior** (modeled by $f_\theta(x, t)[i] \in \Delta(\Sigma)$): for each index i that $x^i = \mathbf{m}$, the posterior
 147 distribution over the underlying clean token.
- 148 • **Insertion expectation** (modeled by $g_\theta(x, t)[i] \in \mathbb{R}_{\geq 0}$): for all indices i in x , the expected number
 149 of tokens that remain to be inserted in between x^{i-1} and x^i .

150 f_θ resembles the familiar unmasking posterior from MDMs, whereas g_θ is new: it predicts how many
 151 tokens need to be inserted. Intuitively, modeling a *variable-length* p_1 is harder than the fixed-length
 152 setup of MDM—introducing an **insertion expectation** allows us to parameterize more complicated
 153 CTMC for FlexMDM; its rate matrix will appear soon in Proposition 2. To define the training loss, we
 154 set the boundary values of s_t as $s_t[-1] := 0$ and $s_t[\text{len}(s_t)] := \text{len}(x_1)$, and let $\phi(x, y) = y - x \log y$
 155 denote a scalar Bregman divergence.

$$\mathcal{L}_\theta = - \int_0^1 \mathbb{E} \left[\underbrace{\frac{\hat{\beta}_t}{1 - \beta_t} \sum_{i: x_t^i = \mathbf{m}} \log f_\theta(x_t, t)[i, x_1^{s_t^i[i]}]}_{\text{unmasking loss}} + \underbrace{\frac{\hat{\alpha}_t}{1 - \alpha_t} \sum_{i=0}^{\text{len}(x_t)} \phi(s_t[i] - s_t[i-1] - 1, g_\theta(x_t, t)[i])}_{\text{insertion loss}} \right] dt. \quad (7)$$

156 Here, the expectation is taken over $x_1 \sim p_1$, $(x_t, s_t) \sim p_t(\cdot | x_1)$. Proposition 1 exactly characterizes
 157 the unmasking posterior and insertion expectation and shows they uniquely minimize (7).

158 **Proposition 1** (FlexMDM training loss). *The loss \mathcal{L}_θ in (7) is uniquely minimized at*

$$f_\theta(x, t)[i, v] = \underbrace{\mathbb{P}(x_1^{s_t^i[i]} = v | x_t = x)}_{\text{unmasking posterior}}, \quad g_\theta(x, t)[i] = \underbrace{\mathbb{E}[s_t[i] - s_t[i-1] - 1 | x_t = x]}_{\text{insertion expectation}}.$$

159 These quantities match the explanation above: the posterior over the clean token together with the
 160 expected number of insertions. They precisely determine the FlexMDM rate matrix stated next.

161 **Proposition 2** (FlexMDM Rate Matrix). *Let the rate matrix R_t be defined as:*

$$\begin{aligned} \text{Unmask} : R_t(x, x[x^i \leftarrow v]) &= \frac{\hat{\beta}_t}{1 - \beta_t} \cdot \mathbb{P}(x_1^{s_t^i[i]} = v | x_t = x), \quad v \in \Sigma, x^i = \mathbf{m} \\ \text{Insert} : R_t(x, x \triangleleft_i \mathbf{m}) &= \frac{\hat{\alpha}_t}{1 - \alpha_t} \cdot \mathbb{E}[s_t[i] - s_t[i-1] - 1 | x_t = x], \end{aligned} \quad (8)$$

162 where $x \triangleleft_i \mathbf{m}$ is the sequence obtained from x by inserting a mask token in between (x^{i-1}, x^i) .
 163 Then R_t solves the KFE (equation (3)) with p_t as the probability mass function of the FlexMDM
 164 interpolant x_t .

165 Proposition 1 thus implies that minimizing the loss yields exact recovery of the rate matrix. In
 166 practice, we could simulate the CTMC using the learned networks (f_θ, g_θ) in place of the ground-
 167 truth quantities in (8). By denoting the resulting terminal distribution as p_1^θ , the variational loss
 168 quantifies the terminal-time KL divergence:

$$\mathcal{D}_{\text{KL}}(p_1 || p_1^\theta) \leq \mathcal{L}_\theta - \mathcal{L}_*$$

169 We defer formal demonstration of propositions and the KL divergence guarantee to Appendix D.
 170 Definition of the joint interpolant is reinstated in definition D.2, the rate matrix in proposition D.3,
 171 the loss and variational bound in proposition D.4.

172 **Remark.** Our FlexMDM interpolant introduces only one extra quantity beyond MDM’s unmasking
 173 posterior: the insertion expectation, a simple scalar per position. This stems from our design choice
 174 to gradually insert and then unmask a token. As shown in Section 5.2, this enables efficient task
 175 transfer of pretrained MDM weights. In contrast, alternative interpolants would require modeling
 176 more complex objects, such as a full token distribution, adding unnecessary training burden.

177 4 Variable Length Masked Diffusions: Inference

178 In this section, we outline inference algorithms for FlexMDM, focusing on two variants: **vanilla**
179 **inference** and **adaptive inference**. We begin with a brief overview of vanilla and adaptive inference
180 in MDMs.

181 **MDM adaptive inference.** For the case of MDM, MDM inference proceeds by simulating the rate
182 matrix entries in 5. From a high-level one way this can be done is by (a) independently sampling
183 a subset of masked tokens to unmask and (b) sampling clean tokens from the unmasking posterior.
184 Crucially for what follows, the same guarantee holds for non-independent *adaptive* choices of
185 unmasking indices, e.g., confidence-based: correctness hinges on using the ground-truth unmasking
186 posterior, not on following the rate matrix’s unmasking entries. This adaptive inference strategy
187 is widely used due to its empirical performance. We adopt this template and show that FlexMDM
188 inherits the same any-order property.

189 **Vanilla inference.** We begin with the *vanilla inference* of FlexMDM, which is obtained by discretizing
190 the CTMC in (8) using trained neural networks (f_θ, g_θ) . Choosing an appropriate discretization
191 scheme is crucial, as different schemes can lead to markedly different empirical behavior. We adopt
192 τ -*leaping*—originating in chemical physics and shown to outperform naive Euler discretization for
193 MDMs [Campbell et al., 2022]—which batches all events occurring within a fixed interval $[t, t + \tau]$.
194 At a high level, for each discretized step, we simultaneously:

- 195 • **Unmasking:** For each mask token, sample for every unmasking a number according to the unmask-
196 ing intensities in the rate matrix. Unmask only if a non-zero entry is returned.
- 197 • **Insertion:** Sample the number of *mask-token insertions* from a Poisson distribution parameterized
198 by the insertion rate, then apply those insertions.

199 As the number of steps $\rightarrow \infty$, this inference algorithm recovers the CTMC and the discretization
200 error vanishes. Algorithm 3 details the full sampler.

201 **Adaptive inference.** Notably, one can choose the positions to unmask **adaptively**. Precisely, at each
202 inference step we select the unmasking positions according to a heuristic rule that prioritizes **the most**
203 **confident indices** (Algorithm 3), where confidence is computed either from the *model’s unmasking*
204 *posterior* or via a *semi-autoregressive* rule (prioritizing leftmost masks). We find such adaptive choice
205 substantially boosts performance; see Section 5.

206 Since unmasking indices in an adaptive no longer trace the transitions described by the rate matrix
207 entries defined in (8), one might ask whether sampling still guarantees to sample from the target
208 distribution p_1 in the infinitesimal limit. The following proposition answers in the affirmative.

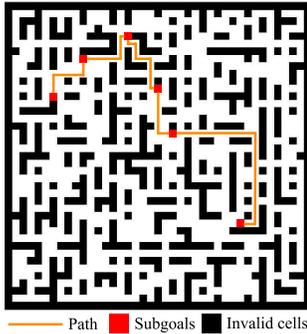
209 **Proposition 3** (Any-order inference, informal). *Consider any sampling scheme that, at each step: (i)*
210 *unmasks an arbitrary subset of masked positions but draws revealed tokens from the ground-truth*
211 *unmasking posterior; and (ii) applies insertion CTMC governed by the ground-truth rate matrix.*
212 *Then the resulting process samples from the target distribution p_1 .*

213 The formal statement and the proof of Proposition 3 are given in Appendix E. In words, following the
214 unmasking entries of the rate matrix corresponding to the schedule used in training is *not* necessary
215 to preserve the sampling guarantee. Moreover, the samplers as $N \rightarrow \infty$ in Algorithm 3 is subsumed
216 by the class in Proposition 3, therefore, assuming access to the ground-truth unmasking posterior and
217 insertion expectation, the corresponding class of algorithms in Algorithm 3 samples from p_1 up to
218 discretization error.

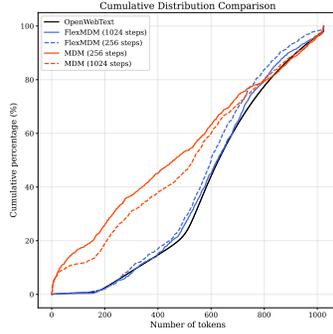
219 **Remark.** A key technical ingredient underlying the rigor of our adaptive inference is that the
220 respective entries of the unmasking posterior of the ground truth rate matrix in Proposition 3 do *not*
221 depend on the choice of unmasking schedule β_t (the proof is given in Appendix E.2.1 and is also
222 shown in the continuous setting in Albergo et al. [2023b], Negrel et al. [2025]). This independence
223 allows a single model f_θ to learn all possible unmasking transitions arising along different paths that
224 ultimately connect p_0 to p_1 , thereby **enabling adaptive unmasking** at inference time. This feature
225 is the same mechanism enabling adaptive inference for MDMs, but for FlexMDMs, proving that it
226 interfaces correctly with insertions is quite subtle. We defer further discussions to Appendix E.

227 5 Experiment

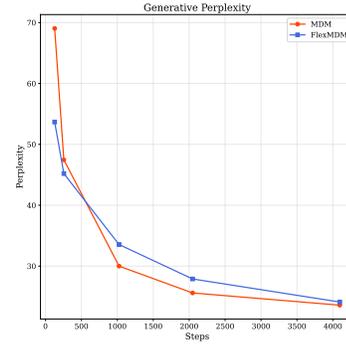
228 In this section, we present experimental results for FlexMDM, demonstrating the following:



(a) Maze task illustration.



(b) FlexMDM recovers the length distribution of pretraining data.



(c) FlexMDM shows comparable generative perplexity to MDM.

Figure 3

- 229 • **FlexMDM is an effective variable-length learner**: length modeling, planning, local edits.
- 230 • **FlexMDM is scalable**: 8B FlexMDM is obtainable by initializing from a pretrained MDM.

231 Section 5.1 presents from-scratch results for FlexMDM on text and planning tasks distributions,
 232 confirming its practical efficiency. Next, Section 5.2 provides an 8B-scale FlexMDM’s training recipe,
 233 initialized from LLaDA-8B Nie et al. [2025], and evaluates it in math and coding tasks. We begin
 234 with the architectural and scheduling choices used throughout.

235 **Training design.** Recall from Section 3 that FlexMDM models the unmasking posterior f_θ and
 236 insertion expectation g_θ given state x and time step t . We adopt DiT [Peebles and Xie, 2023], a
 237 bidirectional transformer that enables additional embedding, as a backbone. To learn both quantities
 238 jointly, we attach two output heads: a standard posterior head for f_θ and a scalar softplus head for g_θ .
 239 Moreover, we choose our unmasking and insertion schedule to be both linear, $\alpha_t = \beta_t = t^2$.

240 5.1 Pretraining

241 In this section, we evaluate FlexMDM’s ability to learn variable-length data from scratch. Our
 242 baseline is MDM, which is fixed-length but can handle variable-length sequences by padding to a
 243 fixed maximum length with an auxiliary pad token. This padding setup is widely used in instruction
 244 fine-tuning when variable-length answers are desired [Nie et al., 2025, 2024, Ye et al., 2025, Gong
 245 et al., 2024]. For a fair comparison, we use vanilla inference for both MDM and FlexMDM throughout.
 246 Further experimental details appear in Appendix F.

247 5.1.1 Pretraining on text data

248 We first construct a training dataset from the raw OpenWebText corpus [Gokaslan et al., 2019], splitting
 249 each article into paragraphs to preserve semantic coherence and yield variable-length sequences.
 250 Models pretrained on this data therefore generate variable-length text.

251 **Results.** We train 175M FlexMDM and MDM with a maximum sequence length 1024 for 500K
 252 iterations and batch size 1024. Using the pretrained models, we vary the number of sampling
 253 steps and measure (a) generative perplexity as a proxy for text fluency, and (b) the induced length distribu-
 254 tion. Figure 3c shows comparative generative perplexity for the two models, improving as sampling
 255 steps increase, indicating no fluency degradation for FlexMDM despite its more involved loss ob-
 256 jective³. Crucially, we observe that **FlexMDM matches the true length distribution far more closely**
 257 (Figure 3b): with only 256 steps it tracks the ground truth distribution, whereas MDM remains
 258 miscalibrated even at 1024 steps.

259 5.1.2 Planning task

260 We further evaluate FlexMDM’s ability in a planning task in a discrete space. Motivated by an
 261 earlier study Janner et al. [2022] that investigated the ability of continuous diffusion in maze tasks,
 262 we design a grid-maze benchmark: the maze is fixed but unknown to the model, with a subset of

²The ground-truth unmasking posterior is independent of β_t , so we condition the network on α_t only; under the linear choice $\alpha_t = t$, this coincides with the usual time embedding.

³We notice that MDM struggles to generate short sentences with low perplexities. Therefore, we filter overly short sequences with ≤ 10 tokens when calculating average perplexity.

263 cells invalid. Given a sequence of subgoal grids (g_1, \dots, g_K) , the model must connect this sequence
 264 without entering invalid cells (see Figure 3a). This subgoal structure aligns naturally with FlexMDM:
 265 starting from (g_1, \dots, g_K) , inference inserts mask tokens between subgoals and then unmask to
 266 generate a feasible path. By contrast, a fixed-length MDM must preassign each subgoal to a specific
 267 position, which is difficult to know *a priori*. We provide additional details in Appendix F.2.

Results. We use a 41×41 maze and control the task’s difficulty via varying the number of subgoals $K \in \{2, 7, 12\}$. As K increases, MDM performance degrades markedly, while FlexMDM maintains robust success rates, reaching a gap of up to 60% at $K = 12$. These results firmly support FlexMDM as a principled approach for subgoal-based planning, where preallocating token positions is inherently challenging for fixed-length models.

Difficulty	MDM	FlexMDM
Easy	68.4%	92.3%
Medium	29.3%	90.4%
Hard	24.2%	90.0%

Table 1: FlexMDM outperforms MDM on the subgoal-style maze-planning task.

269 5.2 Scaling up FlexMDM

270 assess FlexMDM’s scalability by extending it to 8B parameters, observing clear gains over the MDM
 271 baseline. Since both share the unmasking posterior, we hypothesize effective task transfer from a pre-
 272 trained MDM. Starting from LLaDA-Base [Nie et al., 2025], we (a) add time-embedding layers and a
 273 scalar head for insertion expectation, and (b) attach LoRA adapters, yielding $\approx 400M$ trainable param-
 274 eters. To cover natural and mathematical language, we train on a 50:50 mix of OpenWebText [Gokaslan
 275 et al., 2019] and Proof-Pile-2 [Azerbaiyev et al., 2023]. Notably, we observe rapid transfer: within
 276 three days on 16 H100s, the model produces variable-length sentences. We then instruction-fine-tune
 277 (IFT) this base FlexMDM for downstream evaluation (Appendix F)⁴.

278 **Results.** For comparison, we train FlexMDM and LLaDA-Base from the same number of IFT pairs.
 279 For math and code, respectively, we IFT on the GSM8K train (Cobbe et al. [2021]; ≈ 8000 pairs) and
 280 the educational split of opc-sft-stage-2 [Huang et al., 2024] (0.1M pairs), for which IFT-ed models
 281 are evaluated in GSM8K-test and HumanEval-infill [Bavarian et al., 2022] in zero-shot. Sampling
 282 is done by confidence-based sampling with sliding window. Notably, as the number of sampling
 283 steps increases, FlexMDM continues to improve, highlighting its strength in reasoning tasks given
 284 sufficient compute—whereas LLaDA’s performance remains flat. Although in this experiment we use
 285 IFT on task-specific pairs, we expect that, given sufficient compute, training on a much more diverse
 286 instruction–answer pairs will yield a more generalized model.

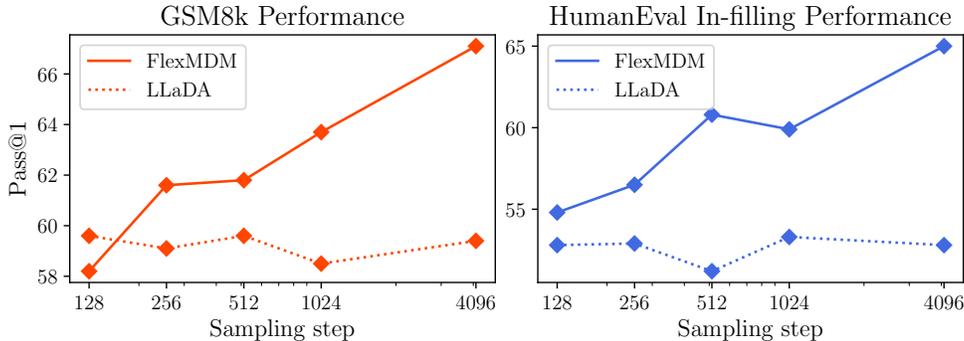


Figure 4: FlexMDM performance exhibits superior scaling when more sampling steps are allocated.

287 6 Conclusion

288 In this work we proposed Flexible Masked Diffusion Models (FlexMDM), a discrete diffusion
 289 framework over variable-length sequences. Theoretically, via a joint interpolant viewpoint, we
 290 provide rigorous guarantees for both training and inference of FlexMDM. Empirically, FlexMDM
 291 learns variable-length structure across diverse scenarios, scales to 8B parameters, trains in only a
 292 few GPU-hours, and yields substantial improvements on math and coding infilling tasks. Further
 293 exploration of FlexMDM’s capabilities is a promising direction for future work.

⁴For fairness, we also IFT LLaDA-Base for the same number of epochs, unlike Zhao et al. [2025].

294 **References**

- 295 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
296 *neural information processing systems*, 33:6840–6851, 2020.
- 297 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
298 Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint*
299 *arXiv:2011.13456*, 2020.
- 300 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised
301 learning using nonequilibrium thermodynamics. In *International conference on machine learning*,
302 pages 2256–2265. pmlr, 2015.
- 303 Jiaxin Shi, Kehang Han, Zhe Wang, Arnaud Doucet, and Michalis Titsias. Simplified and generalized
304 masked diffusion for discrete data. *Advances in neural information processing systems*, 37:
305 103131–103167, 2024.
- 306 Subham Sahoo, Marianne Arriola, Yair Schiff, Aaron Gokaslan, Edgar Marroquin, Justin Chiu,
307 Alexander Rush, and Volodymyr Kuleshov. Simple and effective masked diffusion language
308 models. *Advances in Neural Information Processing Systems*, 37:130136–130184, 2024.
- 309 Itai Gat, Tal Remez, Neta Shaul, Felix Kreuk, Ricky TQ Chen, Gabriel Synnaeve, Yossi Adi, and
310 Yaron Lipman. Discrete flow matching. *Advances in Neural Information Processing Systems*, 37:
311 133345–133385, 2024.
- 312 Jiacheng Ye, Jiahui Gao, Shansan Gong, Lin Zheng, Xin Jiang, Zhenguo Li, and Lingpeng Kong.
313 Beyond autoregression: Discrete diffusion for complex reasoning and planning. *arXiv preprint*
314 *arXiv:2410.14157*, 2024.
- 315 Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-
316 Rong Wen, and Chongxuan Li. Large language diffusion models. *arXiv preprint arXiv:2502.09992*,
317 2025.
- 318 Jiacheng Ye, Zhihui Xie, Lin Zheng, Jiahui Gao, Zirui Wu, Xin Jiang, Zhenguo Li, and Lingpeng
319 Kong. Dream 7b, 2025. URL <https://hkunlp.github.io/blog/2025/dream>.
- 320 Shen Nie, Fengqi Zhu, Chao Du, Tianyu Pang, Qian Liu, Guangtao Zeng, Min Lin, and Chongxuan
321 Li. Scaling up masked diffusion models on text. *arXiv preprint arXiv:2410.18514*, 2024.
- 322 Michael S. Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants,
323 2022.
- 324 Michael S Albergo, Nicholas M Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A unifying
325 framework for flows and diffusions. *arXiv preprint arXiv:2303.08797*, 2023a.
- 326 Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching
327 for generative modeling, 2022.
- 328 Kaiwen Zheng, Yongxin Chen, Hanzi Mao, Ming-Yu Liu, Jun Zhu, and Qinsheng Zhang. Masked
329 diffusion models are secretly time-agnostic masked models and exploit inaccurate categorical
330 sampling. *arXiv preprint arXiv:2409.02908*, 2024.
- 331 Jingyang Ou, Shen Nie, Kaiwen Xue, Fengqi Zhu, Jiacheng Sun, Zhenguo Li, and Chongxuan Li.
332 Your absorbing discrete diffusion secretly models the conditional distributions of clean data. *arXiv*
333 *preprint arXiv:2406.03736*, 2024.
- 334 Andrew Campbell, Joe Benton, Valentin De Bortoli, Thomas Rainforth, George Deligiannidis, and
335 Arnaud Doucet. A continuous time framework for discrete denoising models. *Advances in Neural*
336 *Information Processing Systems*, 35:28266–28279, 2022.
- 337 Michael S. Albergo, Nicholas M. Boffi, Michael Lindsey, and Eric Vanden-Eijnden. Multimarginal
338 generative modeling with stochastic interpolants, 2023b. URL [https://arxiv.org/abs/2310.](https://arxiv.org/abs/2310.03695)
339 03695.

- 340 Hugo Negrel, Florentin Coeurdoux, Michael S. Albergo, and Eric Vanden-Eijnden. Multitask learning
341 with stochastic interpolants, 2025. URL <https://arxiv.org/abs/2508.04605>.
- 342 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*
343 *the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023.
- 344 Shansan Gong, Shivam Agarwal, Yizhe Zhang, Jiacheng Ye, Lin Zheng, Mukai Li, Chenxin An,
345 Peilin Zhao, Wei Bi, Jiawei Han, et al. Scaling diffusion language models via adaptation from
346 autoregressive models. *arXiv preprint arXiv:2410.17891*, 2024.
- 347 Aaron Gokaslan, Vanya Cohen, Ellie Pavlick, and Stefanie Tellex. Openwebtext corpus. [http://](http://Skylion007.github.io/OpenWebTextCorpus)
348 Skylion007.github.io/OpenWebTextCorpus, 2019.
- 349 Michael Janner, Yilun Du, Joshua Tenenbaum, and Sergey Levine. Planning with diffusion for
350 flexible behavior synthesis. In *International Conference on Machine Learning*, 2022.
- 351 Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster, Marco Dos Santos, Stephen McAleer, Albert Q.
352 Jiang, Jia Deng, Stella Biderman, and Sean Welleck. Llemma: An open language model for
353 mathematics, 2023.
- 354 Siyan Zhao, Devaansh Gupta, Qinqing Zheng, and Aditya Grover. d1: Scaling reasoning in diffusion
355 large language models via reinforcement learning. *arXiv preprint arXiv:2504.12216*, 2025.
- 356 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
357 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve
358 math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 359 Siming Huang, Tianhao Cheng, Jason Klein Liu, Jiaran Hao, Liuyihan Song, Yang Xu, J. Yang,
360 J. H. Liu, Chenchen Zhang, Linzheng Chai, Ruifeng Yuan, Zhaoxiang Zhang, Jie Fu, Qian Liu,
361 Ge Zhang, Zili Wang, Yuan Qi, Yinghui Xu, and Wei Chu. Opencoder: The open cookbook for
362 top-tier code large language models. 2024. URL <https://arxiv.org/pdf/2411.04905>.
- 363 Mohammad Bavarian, Heewoo Jun, Nikolas Tezak, John Schulman, Christine McLeavey, Jerry
364 Tworek, and Mark Chen. Efficient training of language models to fill in the middle. *arXiv preprint*
365 *arXiv:2207.14255*, 2022.
- 366 Emiel Hooeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows
367 and multinomial diffusion: Learning categorical distributions. *Advances in neural information*
368 *processing systems*, 34:12454–12465, 2021.
- 369 Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured
370 denoising diffusion models in discrete state-spaces. *Advances in neural information processing*
371 *systems*, 34:17981–17993, 2021.
- 372 Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion modeling by estimating the ratios
373 of the data distribution. *arXiv preprint arXiv:2310.16834*, 2023.
- 374 Andrew Campbell, Jason Yim, Regina Barzilay, Tom Rainforth, and Tommi Jaakkola. Generative
375 flows on discrete state-spaces: Enabling multimodal flows with applications to protein co-design.
376 *arXiv preprint arXiv:2402.04997*, 2024.
- 377 Neta Shaul, Itai Gat, Marton Havasi, Daniel Severo, Anuroop Sriram, Peter Holderrieth, Brian Karrer,
378 Yaron Lipman, and Ricky TQ Chen. Flow matching with general discrete paths: A kinetic-optimal
379 perspective. *arXiv preprint arXiv:2412.03487*, 2024.
- 380 Yuxuan Song, Zheng Zhang, Cheng Luo, Pengyang Gao, Fan Xia, Hao Luo, Zheng Li, Yuehang
381 Yang, Hongli Yu, Xingwei Qu, et al. Seed diffusion: A large-scale diffusion language model with
382 high-speed inference. *arXiv preprint arXiv:2508.02193*, 2025.
- 383 Google DeepMind. Gemini diffusion, 2025. URL [https://blog.google/technology/](https://blog.google/technology/google-deeppmind/gemini-diffusion/)
384 [google-deeppmind/gemini-diffusion/](https://blog.google/technology/google-deeppmind/gemini-diffusion/).
- 385 Inception Labs, Samar Khanna, Siddhant Kharbanda, Shufan Li, Harshit Varma, Eric Wang, Sawyer
386 Birnbaum, Ziyang Luo, Yanis Miraoui, Akash Palrecha, et al. Mercury: Ultra-fast language models
387 based on diffusion. *arXiv preprint arXiv:2506.17298*, 2025.

- 388 Shansan Gong, Ruixiang Zhang, Huangjie Zheng, Jiatao Gu, Navdeep Jaitly, Lingpeng Kong, and
389 Yizhe Zhang. Diffucoder: Understanding and improving masked diffusion models for code
390 generation. *arXiv preprint arXiv:2506.20639*, 2025.
- 391 Alexander Swerdlow, Mihir Prabhudesai, Siddharth Gandhi, Deepak Pathak, and Katerina Fragki-
392 adaki. Unified multimodal discrete diffusion. *arXiv preprint arXiv:2503.20853*, 2025.
- 393 Jaeyeon Kim, Kulin Shah, Vasilis Kontonis, Sham Kakade, and Sitan Chen. Train for the worst, plan
394 for the best: Understanding token ordering in masked diffusions. *arXiv preprint arXiv:2502.06768*,
395 2025.
- 396 Fred Zhangzhi Peng, Zachary Bezemek, Sawan Patel, Jarrid Rector-Brooks, Sherwood Yao,
397 Avishek Joey Bose, Alexander Tong, and Pranam Chatterjee. Path planning for masked diffusion
398 model sampling. *arXiv preprint arXiv:2502.03540*, 2025.
- 399 Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative
400 image transformer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern
401 recognition*, pages 11315–11325, 2022.
- 402 Heli Ben-Hamu, Itai Gat, Daniel Severo, Niklas Nolte, and Brian Karrer. Accelerated sampling from
403 masked diffusion models via entropy bounded unmasking. *arXiv preprint arXiv:2505.24857*, 2025.
- 404 Litu Rout, Constantine Caramanis, and Sanjay Shakkottai. Anchored diffusion language model.
405 *arXiv preprint arXiv:2505.18456*, 2025.
- 406 Long Ma, Fangwei Zhong, and Yizhou Wang. Reinforced context order recovery for adaptive
407 reasoning and planning. *arXiv preprint arXiv:2508.13070*, 2025.
- 408 Zhe Wang, Jiaxin Shi, Nicolas Heess, Arthur Gretton, and Michalis K Titsias. Learning-order autore-
409 gressive models with application to molecular graph generation. *arXiv preprint arXiv:2503.05979*,
410 2025.
- 411 Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and
412 transfer data with rectified flow, 2022. URL <https://arxiv.org/abs/2209.03003>.
- 413 Michael Samuel Albergo, Mark Goldstein, Nicholas Matthew Boffi, Rajesh Ranganath, and Eric
414 Vanden-Eijnden. Stochastic interpolants with data-dependent couplings. In *Forty-first Interna-
415 tional Conference on Machine Learning*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=FFILRGD0jG)
416 [FFILRGD0jG](https://openreview.net/forum?id=FFILRGD0jG).
- 417 Marton Havasi, Brian Karrer, Itai Gat, and Ricky TQ Chen. Edit flows: Flow matching with edit
418 operations. *arXiv preprint arXiv:2506.09018*, 2025.
- 419 Alec Radford and Jeffrey Wu. Rewon child, david luan, dario amodei, and ilya sutskever. 2019.
420 *Language models are unsupervised multitask learners. OpenAI blog*, 1(8):9, 2019.
- 421 Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-
422 efficient exact attention with io-awareness. *Advances in neural information processing systems*, 35:
423 16344–16359, 2022.
- 424 Juechu Dong, Boyuan Feng, Driss Guessous, Yanbo Liang, and Horace He. Flex attention: A
425 programming model for generating optimized attention kernels. *arXiv preprint arXiv:2412.05496*,
426 2024.
- 427 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint
428 arXiv:1711.05101*, 2017.
- 429 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
430 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cris-
431 tian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu,
432 Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,
433 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel
434 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee,
435 Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra,

436 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi,
437 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh
438 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen
439 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic,
440 Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models.
441 *arXiv preprint arXiv: 2307.09288*, 2023.

442 Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
443 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International
444 Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=nZeVKeeFYf9>.
445

446 A Related Works

447 **Discrete diffusion and flows.** Early diffusion models were formulated as continuous-time Markov
448 chains over continuous spaces with Gaussian transition kernels [Sohl-Dickstein et al., 2015, Ho et al.,
449 2020], and were later connected to continuous-time formulations via stochastic differential equations,
450 offering a unifying perspective on score-based generative modeling [Song et al., 2020]. In parallel,
451 *discrete* diffusion has been developed from the viewpoint of Markov chains over discrete space
452 [Hoogetboom et al., 2021]. Notably, Austin et al. [2021] introduced D3PM with several families of
453 discrete transition kernels, and Lou et al. [2023] proposed SEDD, which adopts score-based training
454 objectives. A complementary line of work studies *discrete flows* [Campbell et al., 2024, Gat et al.,
455 2024], aiming to understand continuous-time Markov chains (CTMCs) that interpolate between data
456 and base distributions; this perspective aligns with ours. Subsequent extensions consider token-wise
457 paths and path-wise structure within such flows [Shaul et al., 2024].

458 **Masked Diffusion Models.** Among discrete-transition designs, absorbing-state (a.k.a. masking)
459 kernels have become a popular and strong-performing choice. Recent work shows that this yields
460 a simple and principled training recipe, referred to as Masked Diffusion Models (MDMs) [Sahoo
461 et al., 2024, Shi et al., 2024]. A growing body of results demonstrates the scalability of this approach
462 across problem settings and modalities, including large-scale natural language modeling [Nie et al.,
463 2024, 2025, Ye et al., 2025, Song et al., 2025, DeepMind, 2025], code generation [Labs et al., 2025,
464 Gong et al., 2025], and multimodal learning [Swerdlow et al., 2025].

465 **Any-order inference in MDMs.** With the advent of MDMs, subsequent work has established that
466 they admit theoretically grounded *any-order inference*, wherein tokens can be unmasked in arbitrary
467 orders rather than following a fixed CTMC schedule [Kim et al., 2025, Peng et al., 2025]. Practical
468 token-ordering rules span a spectrum of heuristics based on model confidence and uncertainty—e.g.,
469 maximum-probability logits [Chang et al., 2022, Zheng et al., 2024], probability margin [Kim et al.,
470 2025], semi-autoregressive schedules [Nie et al., 2024], and entropy-based criteria [Ben-Hamu et al.,
471 2025]—as well as strategies that leverage reference models to guide the unmasking trajectory [Peng
472 et al., 2025]. Beyond heuristics, another thread trains auxiliary modules to anchor or adapt the
473 generation order [Rout et al., 2025], while recent work directly *learns* token orders end-to-end [Ma
474 et al., 2025, Wang et al., 2025].

475 **Stochastic interpolant.** Stochastic interpolant Albergo and Vanden-Eijnden [2022], Albergo et al.
476 [2023a] is a general framework for building measure transport based generative models on continuous
477 state space. While building off different philosophical grounds, it can be seen as equivalent to flow
478 matching Lipman et al. [2022], Liu et al. [2022]. Extensions of the interpolant has been proposed for
479 conditional generation through data-dependent coupling Albergo et al. [2024] which we adopt for
480 infilling, and for multi-task learning Negrel et al. [2025] by operator-based formulations.

481 **Descriptive overview on concurrent work.** The most notable concurrent work is EditFlow Havasi
482 et al. [2025] where the primary mathematical machinery that enabled their construction is a novel
483 "Flow Matching with Auxiliary Process". We note that this can be seen to be mathematically
484 equivalent to the notion of joint interpolant in this work. The main difference is philosophical: While
485 EditFlow relies on the notion of probability path as its central mathematical construct, we write down
486 the definition of two coupled random variables which only define the probability path implicitly.

487 Such difference leads to a considerable difference in the interpolation used in the modeling, in the
 488 case of FlexMDM, our choice of the interpolation path enables any-order sampling.

489 B Notation

490 In this section, we reiterate on the notations used in the main body and introduce auxiliary notations
 491 that are used in the proofs of the appendix.

492 **Strings.** Let ε denote the empty string, Σ a vocabulary of words, \mathbf{m} a special mask token. We write
 493 \oplus for string concatenation, e.g., $ab \oplus bc = abbc$, x^i the i -th element of x , $x[S]$ the string indexed by
 494 an index set, e.g. $abc[\{0, 2\}] = ac$. To insert a token v before position i in string x , we write $x \triangleleft_i v$,
 495 e.g., to prepend a token $abc \triangleleft_0 d = dabc$ and to append a token $abc \triangleleft_4 d = abcd$. To replace the i -th
 496 token in x with v , we write $x[x^i \leftarrow v]$. As much of the work involves masking, we write $x \subseteq y$ if x
 497 can be constructed by partially masking y . We write $\text{mask}(x)$, $\text{unmask}(x) \subseteq [\text{len}(x)]$ for the set of
 498 indices corresponding to mask and clean tokens.

499 C Discrete Stochastic Interpolants: Definitions and Propositions for Section 2

500 In continuous spaces, a common approach to define generative transport is the stochastic interpolant
 501 framework, which implicitly defines the interpolation distribution p_t by specifying an interpolant
 502 $\{x_t\}_{t \in [0,1]}$ and regressing the required quantities to realize the transport.

503 In the section C.1, we introduce a discrete analogue of the stochastic interpolant. To illustrate this
 504 framework, we reformulate the widely-used masked diffusion model for sequences of length n within
 505 our setup. Briefly, masked diffusion defines the interpolation p_t by progressively unmasking tokens
 506 in sentences drawn from the data distribution. At time $t = 0$, all tokens are masked, so p_0 is a point
 507 mass at the fully masked sequence the \mathbf{m} token repeated n times. The transition rates driving the
 508 generative transport can be characterized as functions of per-token posterior probabilities conditioned
 509 on time t . These are typically learned by minimizing a variational objective in the form of a weighted
 510 cross-entropy loss.

511 C.1 Discrete Stochastic Interpolant

512 To obtain the target rate matrix, discrete stochastic interpolant relies on an interpolating rate matrix
 513 that drives sample from a sample drawn from a sample from p_0 to a sample from p_1 , defined as
 514 follows:

515 **Definition C.1** (Discrete Stochastic Interpolant and Interpolating Rate). *Let $x_0 \sim p_0$ and $x_1 \sim p_1$.
 516 A **discrete stochastic interpolant** is a family of random variables $\{x_t\}_{t \in [0,1]}$, defined on a common
 517 probability space and satisfying the boundary conditions $x_{t=0} = x_0$ and $x_{t=1} = x_1$, for which there
 518 exists a continuous-time Markov chain with bounded, time-dependent transition rate matrix $K_t^{x_0, x_1}$
 519 such that, for each $t \in [0, 1]$, $\text{Law}(x_t \mid x_0, x_1)$ coincides with the marginal distribution at time t of
 520 that Markov chain started at X_0 . We refer to $K_t^{x_0, x_1}$ as an **interpolating rate matrix**.*

521 With an interpolating rate of an interpolant, the target rate matrix can then be obtained through
 522 Proposition C.1

523 **Proposition C.1** (Target Rate). *Given a discrete stochastic interpolant x_t and an interpolating rate
 524 matrix $K_t^{x_0, x_1}$, the continuous-time Markov chain with initial distribution p_0 and target transition
 525 rate matrix R_t defined as,*

$$R_t(x, y) = \mathbb{E}_{x_0, x_1} [K_t^{x_0, x_1}(x, y) \mid x_t = x]$$

526 *has marginals equal to $\text{Law}(x_t)$.*

527 *Proof.* Writing p_t a pmf of x_t and $p_t(\cdot \mid x_0, x_1)$ a pmf of x_t conditioned on x_0, x_1 . We further write
 528 $q(x_0, x_1)$ the joint pmf of x_0 and x_1 and $q_t(x_0, x_1 \mid x_t)$ to be the joint pmf conditioned on x_t

529 It suffices to show R_t satisfies the Kolmogorov Forward Equation with pmf p_t as follows:

$$\begin{aligned}
\sum_y R_t(y, x) p_t(y) &= \sum_y \mathbb{E}_{x_0, x_1} [K_t^{x_0, x_1}(x, y) | x_t = y] p_t(y) \\
&= \sum_y \sum_{x_0, x_1} K_t^{x_0, x_1}(y, x) q_t(x_0, x_1 | y) p_t(y) \\
&= \sum_y \sum_{x_0, x_1} K_t^{x_0, x_1}(y, x) p_t(y | x_0, x_1) q(x_0, x_1) \\
&= \mathbb{E}_{x_0, x_1} \left[\sum_y K_t^{x_0, x_1}(y, x) p_t(y | x_0, x_1) \right] \\
&= \mathbb{E}_{x_0, x_1} [\partial_t p_t(y | x_0, x_1)] \\
&= \partial_t p_t(x)
\end{aligned}$$

530 This concludes the proof. \square

531 **Remarks.** While written considerably differently, the framework is mathematically equivalent to
532 discrete flow matching. The difference is only philosophical: discrete flow matching relies on the
533 notion of a conditional probability path that the interpolating rate should induce, whereas we define
534 such probability path only implicitly through the definition of the interpolant.

535 C.2 The Masked Diffusion Interpolant

536 As a concrete example of the discrete stochastic interpolant, we reformulate the masked diffusion
537 model and its learning under the framework. As masked diffusion starts from a point mass, we drop
538 the dependence of x_0 in the writing.

539 **Definition C.2** (The Masked Diffusion Interpolant). *Let $x_1 \sim p_1$ be a sentence of length n drawn*
540 *from the data and α_t a smooth unmasking schedule that interpolates from $\alpha_{t=0} = 0$ to $\alpha_{t=1} = 1$.*
541 *Define the unmasking times $\{T^i\}_{i \in [1, \dots, n]}$ as:*

$$\forall i \in \{1, \dots, n\} : T^i \sim \dot{\alpha}_t dt,$$

542 *then, the masked diffusion interpolant is defined as:*

$$x_t = \begin{cases} \mathbf{m} & \text{if } t < T_i, \\ x_1^i & \text{if } t \geq T_i. \end{cases}$$

543 In words, at each time t , x_t reveals a subset of the tokens from x_1 , with each token x_1^i independently
544 unmasked at its associated time T^i .

545 **Proposition C.2** (The Masked Diffusion Interpolating Rate). *One interpolating rate $K_t^{x_1}$ of the*
546 *masked diffusion interpolant x_t is given by:*

$$\forall x \subseteq x_1, v \in \Sigma, x^i = \mathbf{m} : K_t(x, x[x^i \leftarrow v]) = \frac{\dot{\alpha}_t}{1 - \alpha_t} \mathbf{1}\{v = x_1^i\}.$$

547 *Proof.* Let $p_t(\cdot | x_1)$ as the pmf of x_t conditioned on x_1 , from the definition of the interpolant, we
548 notice that:

$$p_t(x | x_1) = \prod_{i=0}^{\text{len}(x_1)-1} (1 - \alpha_t) \mathbf{1}\{x^i = \mathbf{m}\} + \alpha_t \mathbf{1}\{x^i = x_1^i\}$$

549 We verify that K_t satisfies the Kolmogorov Forward Equation (3) under the conditioned pmf as
550 follows,

$$\begin{aligned}
\text{L.H.S} &= \partial_t p_t(x | x_1) = \partial_t \prod_{i=0}^{\text{len}(x_1)-1} (1 - \alpha_t) \mathbf{1}\{x^i = \mathbf{m}\} + \alpha_t \mathbf{1}\{x^i = x_1^i\} \\
&= \sum_{i=0}^{\text{len}(x_1)-1} (-\dot{\alpha}_t \mathbf{1}\{x^i = \mathbf{m}\} + \dot{\alpha}_t \mathbf{1}\{x^i = x_1^i\}) \cdot \prod_{j \neq i} (1 - \alpha_t) \mathbf{1}\{x^j = \mathbf{m}\} + \alpha_t \mathbf{1}\{x^j = x_1^j\}
\end{aligned}$$

551

$$\begin{aligned}
\text{R.H.S} &= \sum_{i=0}^{\text{len}(x_t)-1} \mathbf{1}\{x^i = x_1^i\} \frac{\dot{\alpha}_t}{1 - \alpha_t} p_t(x[x^i \leftarrow \mathbf{m}]|x_1) - \mathbf{1}\{x^i = \mathbf{m}\} \frac{\dot{\alpha}_t}{1 - \alpha_t} p_t(x|x_1) \\
&= \sum_{i=0}^{\text{len}(x_t)-1} \mathbf{1}\{x^i = x_1^i\} \frac{\dot{\alpha}_t}{1 - \alpha_t} (1 - \alpha_t) \prod_{j \neq i} (1 - \alpha_t) \mathbf{1}\{x^j = \mathbf{m}\} + \alpha_t \mathbf{1}\{x^j = x_1^j\} \\
&\quad - \sum_{i=0}^{\text{len}(x_t)-1} \mathbf{1}\{x^i = \mathbf{m}\} \frac{\dot{\alpha}_t}{1 - \alpha_t} (1 - \alpha_t) \prod_j (1 - \alpha_t) \mathbf{1}\{x^j = \mathbf{m}\} + \alpha_t \mathbf{1}\{x^j = x_1^j\} \\
&= \text{L.H.S}
\end{aligned}$$

552 This concludes the proof. \square

553 Note that since each T^i is sampled independently from a continuous distribution, the probability that
554 two unmasking times coincide is zero. Thus, only a single token is unmasked in any infinitesimal
555 transition almost surely.

556 **Proposition 4** (The Masked Diffusion Target Rate). *By proposition C.1, a target rate R_t that induces*
557 *Law(x_t) is:*

$$\forall x \subseteq x_1, v \in \Sigma, x^i = \mathbf{m} : \quad R_t(x, x[x^i \leftarrow v]) = \frac{\dot{\alpha}_t}{1 - \alpha_t} \mathbb{P}(x_1^i = v | x_t = x).$$

558 *Proof.* Following proposition C.2, the result follows from invoking proposition C.1. \square

559 To learn an approximation to R_t , we now parameterize an approximate target rate of the form
560 $\hat{R}_t(x, x[x^i \leftarrow v]) = \frac{\dot{\alpha}_t}{1 - \alpha_t} f_\theta(x, t)[i, v]$ where $f_\theta(x, t)[i, v]$ is a learned approximation to the poste-
561 rior $\mathbb{P}(x_1^i = v | x_t = x)$.

562 The target rate can then be characterized by a variational objective that measures the discrepancy
563 between the true and approximate path measures.

564 **Proposition 5** (Variational Loss for Masked Diffusion). *The loss functional defined as:*

$$L[\hat{R}_t] = \int_0^1 \mathbb{E}_{x_1, x_t} \left[-\frac{\dot{\alpha}_t}{1 - \alpha_t} \sum_{i=1}^n \mathbf{1}\{x_t^i = \mathbf{m}\} \log f_\theta(x_t, t)[i, x_1^i] \right] dt,$$

565 *is uniquely minimized when $\hat{R}_t = R_t$, and is connected to the terminal KL-divergence by:*

$$\mathcal{D}_{\text{KL}}(p_1 || \hat{p}_1) \leq L[\hat{R}_t] - L[R_t],$$

566 *where \hat{p}_1 is the approximate data distribution generated by \hat{R}_t .*

567 *Proof.* Let \mathbb{P} and $\hat{\mathbb{P}}$ be the path measures associated with the continuous-time markov chain of the
568 target rate matrix R_t in proposition C.1 and an approximation through a neural network \hat{R}_t . The
569 variational loss follows by expanding the KL-divergence between the two path measures.

$$\begin{aligned}
\mathcal{D}_{\text{KL}}(\mathbb{P} || \hat{\mathbb{P}}) &= \mathbb{E}_{\mathbb{P}} \left[\int_{t=0}^{t=1} R_t(x_t, x_t) - \hat{R}_t(x_t, x_t) dt + \sum_{t: x_t \neq x_{t-}} \log \frac{R_t(x_{t-}, x_t)}{\hat{R}_t(x_{t-}, x_t)} \right] \\
&= \int_0^{t=1} \mathbb{E}_{x_t \sim \mathbb{P}_t} \left[\sum_{y \neq x_t} \hat{R}_t(x_t, y) - R_t(x_t, y) \log \hat{R}_t(x_t, y) \right] dt + \text{Const.} \\
&= \int_0^1 \mathbb{E}_{x_1, x_t} \left[-\frac{\dot{\alpha}_t}{1 - \alpha_t} \sum_{i=1}^n \mathbf{1}\{x_t^i = \mathbf{m}\} \log f_\theta(x_t, t)[i, x_1^i] \right] dt + \text{Const.}
\end{aligned}$$

570 where the first line takes expectation of the Radon-Nikodym derivative between the two path measures.

571 The terminal KL bound then follows directly from the data processing inequality. \square

572 **D Joint Discrete Stochastic Interpolants: Definitions and Propositions for**
573 **Section 3**

574 Building on the discrete stochastic interpolant, we proceed to construct a masked diffusion model
575 which natively supports growing in length by inserting mask tokens in arbitrary positions.

576 On a high-level, we would like define an interpolant constructed by deleting and masking sentences
577 from the data distribution. However, the corresponding interpolating rate becomes hard to define, as
578 it is no longer clear what each mask token should unmask to.

579 To this end, we introduce the joint interpolant mechanism that allows us to construct a broader class
580 of interpolant and interpolating rate matrix by augmenting the interpolant with auxiliary information
581 that allows us to specify more flexible interpolation path. We then leverage this newfound freedom to
582 construct the flexible-length masked diffusion model.

583 **D.1 Joint Interpolant**

584 By introducing an auxiliary variable coupled with the interpolant, joint interpolant expands the class
585 of interpolating rate that can be defined.

586 **Definition D.1** (Joint Interpolant and Joint Interpolating Rate). *Let $x_0 \sim p_0$ and $x_1 \sim p_1$. A **joint***
587 ***interpolant** is a family of coupled random variables $\{(x_t, s_t)\}_{t \in [0,1]}$ defined on a common probability*
588 *space and satisfying the boundary conditions $x_{t=0} = x_0$ and $x_{t=1} = x_1$, for which there exists a*
589 *continuous-time Markov chain with bounded, time-dependent transition rate matrix $K_t^{x_0, x_1}$ on the*
590 *joint state space such that, for each $t \in [0, 1]$, the conditional law $\text{Law}(x_t, s_t \mid x_0, x_1)$ coincides*
591 *with the marginal distribution at time t of this Markov chain started at (x_0, s_0) . We call $K_t^{x_0, x_1}$ a*
592 ***joint interpolating rate matrix**.*

593 **Proposition D.1** (Joint Interpolant Target Rate). *Let $\{(x_t, s_t)\}_{t \in [0,1]}$ be a joint interpolant with joint*
594 *interpolating rate matrix $K_t^{x_0, x_1}$. Consider the continuous-time Markov chain with initial distribution*
595 *p_0 and target transition rate matrix R_t defined by*

$$R_t(x, y) = \mathbb{E}_{s_t, x_0, x_1} \left[\sum_{s' \in \mathcal{S}} K_t^{x_0, x_1}((x, s_t), (y, s')) \mid x_t = x \right],$$

596 where \mathcal{S} denotes the discrete state space of the auxiliary variable s_t . The marginal of the chain at
597 time t is then equal to $\text{Law}(x_t)$.

598 *Proof.* Let $p_t(\cdot, \cdot \mid x_0, x_1)$ the joint pmf of x_t, s_t conditioned on x_0, x_1 , let $q_t(x_0, x_1, s_t \mid x_t)$ to be the
599 joint of x_0, x_1, s_t conditioned on x_t , and let $q_t(x_0, x_1)$ to be the joint of x_0, x_1 . We proceed to verify
600 R_t satisfies the KFE as in Equation 3,

$$\begin{aligned} \sum_y R_t(y, x) p_t(y) &= \sum_y \mathbb{E}_{s_t, x_0, x_1} \left[\sum_{s' \in \mathcal{S}} K_t^{x_0, x_1}((y, s_t), (x, s')) \mid x_t = y \right] p_t(y) \\ &= \sum_y \sum_{s_t, x_0, x_1} \sum_{s' \in \mathcal{S}} K_t^{x_0, x_1}((y, s_t), (x, s')) q_t(x_0, x_1, s_t \mid y) p_t(y) \\ &= \sum_y \sum_{s_t, x_0, x_1} \sum_{s' \in \mathcal{S}} K_t^{x_0, x_1}((y, s_t), (x, s')) q_t(x_0, x_1, s_t \mid y) p_t(y) \\ &= \sum_y \sum_{s_t, x_0, x_1} \sum_{s' \in \mathcal{S}} K_t^{x_0, x_1}((y, s_t), (x, s')) p_t(y, s_t \mid x_0, x_1) q(x_0, x_1) \\ &= \mathbb{E}_{x_0, x_1} \left[\sum_{s'} \partial_t p_t(x, s' \mid x_0, x_1) \right] \\ &= \partial_t p_t(x \mid x_0, x_1) \end{aligned}$$

601 This concludes the proof. □

602 **D.2 Flexible-Length Masked Diffusion**

603 We then instantiate the joint interpolant to obtain the length-aware masked diffusion model, using a
 604 sorted list of indices that has been inserted. Again, we drop x_0 in the writing as the model interpolates
 605 between a point mass at an empty sentence to the full data distribution.

606 **Definition D.2** (Flexible-Length Masked Diffusion Joint Interpolant). *Let $x_1 = (x_1^0, \dots, x_1^{n-1}) \sim p_1$
 607 be a sequence of length n . Let α_t and β_t be differentiable schedules on $[0, 1]$ such that $\alpha_0 = \beta_0 = 0$
 608 and $\alpha_1 = \beta_1 = 1$. Define insertion and unmasking times $\{T_1^i\}_{i=1}^n, \{T_2^i\}_{i=1}^n$ as follows:*

$$T_1^i \sim \dot{\alpha}_t dt,$$

$$T_2^i \sim \mathbf{1}\{t \geq T_1^i\} \cdot \frac{\dot{\beta}_t}{1 - \beta_{T_1^i}} dt.$$

609 *At each time $t \in [0, 1]$, define the sorted index set s_t as:*

$$s_t = \{i \in \{1, \dots, n\} \mid t > T_1^i\},$$

610 *with ascending order $s_t[0] < \dots < s_t[|s_t| - 1]$, and boundary values:*

$$s_t[-1] = 0, \quad s_t[|s_t|] = n,$$

611 *and define the interpolant state x_t per-coordinate as:*

$$x_t^i = \begin{cases} \mathbf{m} & \text{if } t < T_2^{s_t[i]}, \\ x_1^{s_t[i]} & \text{if } t \geq T_2^{s_t[i]}. \end{cases}$$

612 *The process $(s_t, x_t)_{t \in [0,1]}$ is the **flexible length masked diffusion joint interpolant**.*

613 Here, s_t tracks the ordered set of indices whose tokens have been inserted and $s_t[i]$ is the i -th smallest
 614 element in s_t . The interpolant x_t reveals the true token x_1^i only after both insertion and unmasking.

615 Given access to this ordered set, one interpolating rate is:

616 **Proposition D.2** (Length-Aware Masked Diffusion Interpolating Rate). *A joint interpolating rate
 617 matrix $Q_t^{x_0, x_1}$ for the joint interpolant above is given by:*

618 *1. Unmask: For index set $s, x \subseteq x_1[s]$, $x^i = \mathbf{m}$, and $v \in \Sigma$:*

$$Q_t^{x_1}((x, s), (x[x^i \leftarrow v], s)) = \frac{\dot{\beta}_t}{1 - \beta_t} \cdot \mathbf{1}\{x_1^{s[i]} = v\}$$

619 *2. Insert: For index set $s, x \subseteq x_1[s]$, $j \notin s$, and position i such that $s[i-1] < j < s[i]$:*

$$Q_t^{x_1}((x, s), (x \triangleleft_i \mathbf{m}, s \cup \{j\})) = \frac{\dot{\alpha}_t}{1 - \alpha_t}$$

620 *Proof.* We first write down the $p_t(\cdot, \cdot | x_1)$ the conditioned pmf of (s, x) given x_0, x_1 . Let $n = \text{len}(x_t)$,
 621 then

$$p_t(s, x | x_1) = A(t) \prod_{i \in s_t} I_i(t),$$

$$A(t) := (1 - \alpha_t)^{n - |s|}$$

$$I_i(t) := \int_0^t \dot{\alpha}_u \left(\frac{1 - \beta_t}{1 - \beta_u} \mathbf{1}\{x^i = \mathbf{m}\} + \frac{\beta_t - \beta_u}{1 - \beta_u} \mathbf{1}\{x^i = x_1^i\} \right) du.$$

622 Differentiate using the product rule:

$$\partial_t p_t(s, x | x_1) = \dot{A}(t) \prod_{i \in s} I_i(t) + A(t) \sum_{j \in s} \left(\dot{I}_j(t) \prod_{i \in s \setminus \{j\}} I_i(t) \right).$$

623 For $A(t) = (1 - \alpha_t)^{n-|s|}$ we have

$$\dot{A}(t) = -(n - |s|)\dot{\alpha}_t(1 - \alpha_t)^{n-|s|-1} = A(t)\left(-\frac{(n - |s|)\dot{\alpha}_t}{1 - \alpha_t}\right).$$

624 For each $i \in s$ apply the Leibniz rule to $I_i(t)$:

$$\begin{aligned} \dot{I}_i(t) &= \dot{\alpha}_t \mathbf{1}\{x^i = \mathbf{m}\} + \int_0^t \dot{\alpha}_u \left(\frac{-\dot{\beta}_t}{1 - \beta_u} \mathbf{1}\{x^i = \mathbf{m}\} + \frac{\dot{\beta}_t}{1 - \beta_u} \mathbf{1}\{x^i = x_1^i\} \right) du \\ &= \dot{\alpha}_t \mathbf{1}\{x^i = \mathbf{m}\} + \dot{\beta}_t \left(-\mathbf{1}\{x^i = \mathbf{m}\} + \mathbf{1}\{x^i = x_1^i\} \right) \int_0^t \frac{\dot{\alpha}_u}{1 - \beta_u} du. \end{aligned}$$

625 Substituting $\dot{A}(t)$ and $\dot{I}_i(t)$ into the product-rule expansion yields

$$\begin{aligned} \partial_t p_t(s, x | x_1) &= -(n - |s|)\dot{\alpha}_t(1 - \alpha_t)^{n-|s|-1} \prod_{i \in s} I_i(t) \\ &\quad + (1 - \alpha_t)^{n-|s|} \sum_{j \in s} \left[\left(\dot{\alpha}_t \mathbf{1}\{x^j = \mathbf{m}\} + \dot{\beta}_t (-\mathbf{1}\{x^j = \mathbf{m}\} + \mathbf{1}\{x^j = x_1^j\}) \int_0^t \frac{\dot{\alpha}_u}{1 - \beta_u} du \right) \right. \\ &\quad \left. \cdot \prod_{i \in s \setminus \{j\}} I_i(t) \right]. \end{aligned}$$

626 This can then be rewritten as

$$\begin{aligned} \partial_t p_t(s, x | x_1) &= - \sum_{i=0}^{\text{len}(x_t)} \frac{\dot{\alpha}_t}{1 - \alpha_t} p_t(s, x | x_1) - \sum_{x^i = \mathbf{m}} \frac{\dot{\beta}_t}{1 - \beta_t} p_t(s, x | x_1) \\ &\quad + \sum_{x^i \neq \mathbf{m}} \frac{\dot{\beta}_t}{1 - \beta_t} p_t(s, x[x^i \leftarrow \mathbf{m}] | x_1) + \sum_{x^i = \mathbf{m}} \frac{\dot{\alpha}_t}{1 - \alpha_t} p_t(s - \{s[i]\}, \text{remove}(x, i) | x_1) \end{aligned}$$

627 where $\text{remove}(x, i)$ refers to the string constructed by removing the i -th element of x .

628 Notice that this is equivalent to the R.H.S of the KFE (Eq. 3) if one uses the rate matrix R_t . This
629 concludes the proof.

630 □

631 **Proposition D.3** (FlexMDM Rate Matrix (Restated from Proposition 2)). *By Proposition D.1, the*
632 *induced marginal target rate R_t is:*

633 1. *Unmask: For $x^i = \mathbf{m}$, $v \in \Sigma$:*

$$R_t(x, x[x^i \leftarrow v]) = \frac{\dot{\beta}_t}{1 - \beta_t} \cdot \mathbb{P}(x_1^{s_t[i]} = v | x_t = x)$$

634 2. *Insert: For position $i \in \{0, \dots, |x|\}$:*

$$\bar{R}_t(x, x \triangleleft_i \mathbf{m}) = \frac{\dot{\alpha}_t}{1 - \alpha_t} \cdot \mathbb{E}_{s_t} [s_t[i] - s_t[i - 1] - 1 | x_t = x]$$

635 *Proof.* The proof follows by noting proposition D.2 and invoking proposition D.1. □

636 Parametrizing a approximate target rate matrix \hat{R}_t in terms of an approximate per-token posterior
637 $f_\theta(x, t)[i, v] \approx \mathbb{P}(x_1^{s_t[i]} = v | x_t = x)$, and an approximate number of insertion $g_\theta(x, t)[i] \approx$
638 $\mathbb{E}_{s_t} [s_t[i] - s_t[i - 1] - 1 | x_t = x]$. The target rate matrix can be learnt by minimising the following
639 variational objective.

640 **Proposition D.4** (FlexMDM Loss (Restated from Proposition 1)). *The loss functional defined as:*

$$L[\hat{R}_t] = \int_0^1 \mathbb{E}_{x_1, s_t, x_t} \left[-\frac{\dot{\beta}_t}{1 - \beta_t} \sum_{i=1}^{|x_1|} \mathbf{1}\{x_t^i = \mathbf{m}\} \log f_\theta(x_t, t)[i, v] \right] dt, \\ + \int_0^1 \mathbb{E}_{x_1, s_t, x_t} \left[\frac{\dot{\alpha}_t}{1 - \alpha_t} \sum_{i=0}^{|x_1|} \phi(s_t[i] - s_t[i-1] - 1, g_\theta(x_t, t)[i]) \right] dt,$$

641 where $\phi(x, y) = y - x \log y$, is uniquely minimized when $\hat{R}_t = R_t$ and is connected by terminal
642 KL-divergence by:

$$\mathcal{D}_{\text{KL}}(p_1 || \hat{p}_1) \leq L[\hat{R}_t] - L[R_t],$$

643 where \hat{p}_1 is the approximate data distribution induced by \hat{R}_t .

644 *Proof.* Let \mathbb{P} and $\hat{\mathbb{P}}$ be the path measure of a continuous time markov chain starting with the empty
645 string with rate matrix R_t and \hat{R}_t respectively.

646 Consider the KL-divergence between path measures \mathbb{P} and $\hat{\mathbb{P}}$,

$$\mathcal{D}_{\text{KL}}(\mathbb{P} || \hat{\mathbb{P}}) = \mathbb{E}_{\mathbb{P}} \left[\int_{t=0}^{t=1} R_t(x_t, x_t) - \hat{R}_t(x_t, x_t) dt + \sum_{t: x_t \neq x_{t-}} \log \frac{R_t(x_{t-}, x_t)}{\hat{R}_t(x_{t-}, x_t)} \right] \\ = \int_0^{t=1} \mathbb{E}_{x_t} \left[\sum_{y \neq x_t} \hat{R}_t(x_t, y) - R_t(x_t, y) \log \hat{R}_t(x_t, y) \right] dt + \text{Const.} \\ = \int_0^{t=1} \mathbb{E}_{x_t} \left[\sum_{y \neq x_t} \hat{R}_t(x_t, y) - R_t(x_t, y) \log \hat{R}_t(x_t, y) \right] dt + \text{Const.} \\ = \int_0^1 \mathbb{E}_{x_1, s_t, x_t} \left[-\frac{\dot{\beta}_t}{1 - \beta_t} \sum_{i=1}^{|x_1|} \mathbf{1}\{x_t^i = \mathbf{m}\} \log f_\theta(x_t, t)[i, v] \right] dt \\ + \int_0^1 \mathbb{E}_{x_1, s_t, x_t} \left[\frac{\dot{\alpha}_t}{1 - \alpha_t} \sum_{i=0}^{|x_1|} \phi(s_t[i] - s_t[i-1] - 1, g_\theta(x_t, t)[i]) \right] dt + \text{Const.}$$

647 The terminal KL-bound then follows from the data processing inequality. \square

648 E Details for Section 4

649 E.1 Precise detail on the inference algorithms

650 In this section, we provide details on the unmasking steps of vanilla and adaptive inference for
651 FlexMDM, summarized in Algorithm 3. Suppose at inference time we are given the discretization
652 step size τ , a partially observed sequence x_{t_k} , and the current time step t_k .

653 **Vanilla inference.** For each masked position i (i.e., $x_{t_k}^i = \mathbf{m}$) and each clean token $v \in \Sigma$, we
654 sample unmasking events from a Poisson distribution $\text{Poisson}(R_v \tau)$, where R_v is the unmasking
655 rate toward token v . Concretely, $R_v = \frac{\dot{\beta}_{t_k}}{1 - \beta_{t_k}} \cdot f_\theta(x_{t_k}, t_k)[i, v]$, so that the event count is distributed
656 as $k_v \sim \text{Poi} \left(\tau \cdot \frac{\dot{\beta}_{t_k}}{1 - \beta_{t_k}} \cdot f_\theta(x_{t_k}, t_k)[i, v] \right)$. A masked position is unmasked only if *exactly one* token
657 v produces a count $k_v = 1$ while all others produce zero. This tau-leaping scheme batches all events
658 that occur within the interval $[t_k, t_k + \tau]$.

Algorithm 1 VLMDM inference

Require: Learned functions (f_θ, g_θ)
Require: Discretization $0 = t_1 < \dots < t_N = 1$
Require: Insertion, Unmasking schedule α_t, β_t

- 1: Initialize $X_{t_1} \leftarrow \varepsilon$
- 2: **for** $k = 1$ to $N - 1$
- 3: $\tau \leftarrow t_{k+1} - t_k$
- 4: **Invoke Algorithm 2 for unmasking**
- 5: **for** i in $[\text{len}(X_{t_k})]$
- 6: Set rate $r \leftarrow \frac{\dot{\alpha}_{t_k}}{1 - \alpha_{t_k}} \cdot \tau$
- 7: Sample $\ell \sim \text{Poi}(r \cdot g_\theta(X_{t_k}, t_k)[i])$
- 8: **Insert ℓ masks between $X_{t_k}^{i-1}$ and $X_{t_k}^i$**
- 9: **return** X_{t_N}

Algorithm 2 Unmasking Step

- 1: **if vanilla inference :**
- 2: **for** $i \in \{i | X_{t_k}^i = \mathbf{m}\}$ and $v \in \Sigma$
- 3: Set rate $r \leftarrow \frac{\dot{\beta}_{t_k}}{1 - \beta_{t_k}} \cdot \tau$
- 4: $k_v \sim \text{Poi}(r \cdot f_\theta(X_{t_k}, t_k)[i, v])$
- 5: **if** $\exists! v$ such that $k_v = 1$
- 6: Set $X_{t_k}^i \leftarrow v$
- 7: **if adaptive inference :**
- 8: Select K (the size of $|S|$)
- 9: **for** $i \in \{i | X_{t_k}^i = \mathbf{m}\}$
- 10: Compute confidence \mathcal{C}^i
- 11: **for** i in $\text{argmaxK}(C)$
- 12: $X_{t_k}^i \sim \text{Cat}(f_\theta(X_{t_k}, t_k)[i])$

Algorithm 3: VLMDM inference. At each step we perform **unmasking** and **insertion**. For **unmasking**, unmask by τ -leaping (**vanilla**) or by confidence-based selection (**adaptive**). The number of mask tokens to **insert** is drawn from a Poisson distribution. **Notation:** Cat , Poi imply the categorical and Poisson distribution, respectively. $\text{argmaxK}(C)$ is the indices set of the K largest components of C .

659 **Adaptive inference.** We first draw the number of tokens to unmask, denoted by an integer K . While
660 Proposition 3 and Theorem E.1 show that the choice of K does not affect the theoretical guarantees,
661 in practice, we set K to match the expected number of unmasked tokens under vanilla inference
662 yields stable behavior. Accordingly, we sample $K \sim \text{Poi}\left(\tau \cdot \frac{\dot{\beta}_{t_k}}{1 - \beta_{t_k}} \cdot \#\{\text{masked tokens in } x_{t_k}\}\right)$.

663 Next, we compute a confidence score for each masked position, based on heuristics such as:

- 664 • Top-K probability [Chang et al., 2022, Zheng et al., 2024]: For state x at time t , the confidence at
665 position i is given by $\max_{v \in \Sigma} f_\theta(x, t)[i, v]$.
- 666 • Top-K probability with sliding window: We further restrict sampling to the leftmost L tokens,
667 where

$$L = \min(\lfloor \gamma_1 \cdot L \rfloor, \gamma_2),$$

668 with γ_1 and γ_2 hyperparameters. This approach is related to semi-autoregressive strategies used in
669 Nie et al. [2025].

670 Finally, we select the subset of positions to unmask as the top- K masked indices with the highest
671 confidence scores.

672 E.2 Proof of FlexMDM’s any-order inference capability.

673 E.2.1 Proof preliminaries

674 **Form of posterior.** We first compute, for each x_t^i , the probabilities of being masked or deleted in
675 equation (6). This follows from a straightforward calculation using the joint distribution of (T_1^i, T_2^i) :

$$p(x_t^i = (\text{empty})) = p(T_1^i > t) = 1 - \alpha_t,$$

$$p(x_t^i = \mathbf{m}) = p(T_1^i \leq t, T_2^i > t) = \int_t^1 \int_0^t \left(\frac{\dot{\beta}_s}{1 - \beta_u} \times \dot{\alpha}_s \right) ds du =: 1 - \gamma_t.$$

676 Here we define γ_t as $1 - p(x_t^i = \mathbf{m})$. Therefore, the process in equation (6) is equivalent to observing
677 a partially masked subsequence x_t obtained by sampling $x_1 \sim p$ and, for each position of x_1 ,
678 independently deleting it with probability $1 - \alpha_t$, masking it with probability $1 - \gamma_t$, or leaving it
679 unchanged with probability $\alpha_t + \gamma_t - 1$. Note that α_t and γ_t both increase from 0 to 1 as t increases
680 from 0 to 1.

681 The posterior is given by

$$\begin{aligned} p(x_1 = x^* \mid x_t = x) & \\ & \propto p(x^*) \cdot p(x_t = x \mid x_1 = x^*) \\ & = p(x^*) \cdot (1 - \alpha_t)^{\text{len}(x^*) - \text{len}(x)} (1 - \gamma_t)^{\#\text{mask}(x)} (\alpha_t + \gamma_t - 1)^{\#\text{unmask}(x)} \cdot \#\{s : x \subseteq x^* \mid s\} \\ & \propto p(x^*) \cdot (1 - \alpha_t)^{\text{len}(x^*) - \text{len}(x)} \cdot \#\{s : x \subseteq x^* \mid s\}. \end{aligned}$$

682 Importantly, as in the vanilla MDM setting, the posterior does not depend on the unmasking schedule
683 (γ_t) (thus β_t), which will enable us to perform unmasking in adaptively chosen positions. Note
684 also that if all sequences in the support of p were of the same length, this posterior would also be
685 independent of (α_t) ; while we do not prove it, in this case this would allow us to choose an arbitrary
686 order of unmaskings and insertions.

687 **Extension of posterior to $t = 1$.** Motivated by the form of the posterior above, we define the
688 following:

$$q_t(x^* | x) \propto \begin{cases} p(x^*) \cdot \mathbf{1}_{x \subseteq x^*} & \text{if } t = 1 \\ p(x^*) \cdot (1 - \alpha_t)^{\text{len}(x^*) - \text{len}(x)} \cdot \#\{s : x \subseteq x^* | s\} & \text{otherwise} \end{cases}$$

689 Note that for $t < 1$, this is the same as $p(x_1 = x | x_t = x)$. We will denote the marginals of $q_t(\cdot | x)$
690 by $q_t^i(\cdot | x)$ for $v \in \Sigma$. The reason for extending the definition of the posterior to $t = 1$ is that in an
691 adaptive FlexMDM sampler (see Definition E.1), because we are entirely decoupling unmasking from
692 the schedule of insertions, after the final insertion step the time parameter t may be 1 even though
693 there are still tokens left to unmask. We will assume oracle access to $q_1(\cdot | x)$ as in practice these are
694 simply the any-order marginals for p , and furthermore in practice these are already well-approximated
695 by the learned posterior marginals $p(x_1^i \cdot | x_{1-\delta} = x)$ for arbitrarily small $\delta > 0$.

696 **Index-tracking variable.** Recall that in the definition of the joint interpolant we defined an index-
697 tracking variable s_t which essentially tracked which indices of x_1 correspond to the tokens in x_t .
698 While our analysis below will not use the language of stochastic interpolants, we will still use the
699 idea of tracking s_t , with slightly modified notation. Specifically, for any $0 \leq t < 1$, we will use
700 the notation $\Pr_{(x_1, s) | x_t = x}$ and $\mathbb{E}_{(x_1, s) | x_t = x}$ to denote probability and expectation with respect to the
701 distribution given by sampling x_1 from $q_t(\cdot | x)$, and then sampling s uniformly random from subsets
702 for which $x \subseteq (x_1)_s$. When we only care about the marginal distribution over s , we write $\Pr_{s | x_t = x}$
703 and $\mathbb{E}_{s | x_t = x}$. Given such a subset s and $i \in \{1, \dots, |s|\}$, we use s_i to denote its i -th element in sorted
704 order. The insertion expectations which we had denoted by $[s_t[i] - s_t[i - 1] - 1]$ in the main body
705 are thus given by $\mathbb{E}_{s | x_t = x} [s_i - s_{i-1} - 1]$ in the notation of this section.

706 E.3 Formal guarantee for adaptive inference

707 **Definition E.1.** Given query access to the posterior marginals $q_t^i(\cdot | x_t = x)$ and to the insertion
708 expectations $\mathbb{E}_{s | x_t = x} [s_i - s_{i-1} - 1]$, an adaptive FlexMDM sampler is any algorithm which produces
709 a sequence of iterates $\hat{x}_{t_1}, \dots, \hat{x}_{t_n}$, where $0 = t_1 < \dots < t_n = 1$, by starting at $\hat{x}_{t_1} = \varepsilon$ and
710 arbitrarily alternating between steps of the following form:

711 • Any-order unmasking step: Starting from \hat{x}_{t_j} , if $\text{mask}(\hat{x}_{t_j})$ is nonempty, pick an arbitrary index i
712 therein (possibly probabilistically), sample v from $q_{t_j}^i(\cdot | \hat{x}_{t_j})$, and set $\hat{x}_{t_j} \leftarrow \hat{x}_{t_j}[\hat{x}_{t_j}^i \leftarrow v]$.

713 • Insertion step: Starting from \hat{x}_{t_j} , run the CTMC with rate matrix

$$R_t^{\text{ins}}(x, y) = \begin{cases} \mathbb{E}_{s | x_t = x} [s_{i+1} - s_i - 1] \cdot \frac{\dot{\alpha}_t}{1 - \alpha_t} & \text{if } y = x \triangleleft_i \mathbf{m} \\ - \sum_{i=0}^{\text{len}(x)} R_t^{\text{ins}}(x, x \triangleleft_i \mathbf{m}) & \text{if } y = x \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

714 for $t_j \leq t \leq t_{j+1}$ to obtain $\hat{x}_{t_{j+1}}$. If $t_{j+1} = 1$, then apply any-order unmasking until $\text{mask}(\hat{x}_{t_{j+1}})$
715 is empty, and terminate.

716 Note that the rate matrix in the second bullet point above is identical to the one in the main body
717 except that we only consider transitions given by insertions.

718 Formally, we will show the following:

719 **Theorem E.1.** Any adaptive FlexMDM sampler for p will generate a sequence of iterates $\hat{x}_{t_1}, \dots, \hat{x}_{t_n}$
720 such that \hat{x}_{t_n} is exactly a sample from p .

721 **E.4 Proof of Theorem E.1**

722 To show that adaptive sampling works, we inductively prove an even stronger statement: at any
 723 intermediate step of the sampler after it has produced \hat{x}_{t_j} , the final output \hat{x}_{t_n} is a sample from
 724 $q_{t_j}(\cdot \mid \hat{x}_{t_j})$.

725 The following two lemmas provide the inductive steps for unmasking and insertion respectively:

726 **Lemma E.1** (Inductive step for unmasking). *Let $0 \leq t \leq 1$ and let x be a partially unmasked*
 727 *subsequence of length m . Let $\pi = (\pi_i)_{i \in \text{mask}(x)}$ denote any distribution over masked indices of x .*
 728 *Suppose that one runs the following:*

- 729 1. Sample index i from π
- 730 2. Sample v from the posterior marginal $q_t^i(\cdot \mid x_t = x)$
- 731 3. Sample from $q_t(\cdot \mid x[x^i \leftarrow v])$.

732 *The output of this procedure is a sample x^* from $q_t(\cdot \mid x)$.*

733 **Lemma E.2** (Inductive step for insertion). *Let $0 \leq t < 1$ and let x be a partially unmasked*
 734 *subsequence of length m . Let $0 \leq h \leq 1 - t$ be any duration of time. Suppose that one runs the*
 735 *following:*

- 736 1. Starting from x , run the CTMC with rate matrix given by Eq. (9) for time h to obtain x'
- 737 2. Sample from $q_{t+h}(\cdot \mid x')$.

738 *The output of this procedure is a sample from the posterior $q_t(\cdot \mid x)$.*

739 We defer the proofs of these to Sections E.5 and E.6 below. The idea for the former is identical to
 740 the proof of the folklore fact that vanilla MDMs can sample in any order [Kim et al., 2025]. The
 741 proof for the latter is more involved and involves explicitly verifying that the Kolmogorov *backward*
 742 equation is satisfied by the rate matrix we have constructed.

743 Here we verify that these Lemmas are enough to establish Theorem E.1.

744 *Proof of Theorem E.1.* We show more generally that starting from any intermediate time step \hat{x}_{t_j}
 745 (not just $j = 1$), any adaptive FlexMDM sampler outputs a sample from $q_{t_j}(\cdot \mid \hat{x}_{t_j})$. We do this by
 746 inducting on the total number of insertion steps that remain.

747 As the base case for the induction, if no more insertion steps remain, then we must have $t_j = 1$. In
 748 this case, we can further induct on the number of unmasking steps and apply Lemma E.1 with t
 749 therein set to 1 to conclude that the final output is a sample from $q_1(\cdot \mid \hat{x}_{t_j})$.

750 For the inductive step, we have $t_j < 1$ and suppose we have shown that for any FlexMDM sampler
 751 that makes at most m insertion steps, starting from any \hat{x}_{t_j} at intermediate time t_j , it samples from
 752 $q_{t_j}(\cdot \mid \hat{x}_{t_j})$. Now consider a FlexMDM sampler that makes at most $m + 1$ insertion steps starting
 753 from \hat{x}_{t_j} at intermediate time t_j . If in the next step it performs an insertion step, i.e. it runs the
 754 CTMC with rate matrix defined above for total time $h = t_{j+1} - t_j$, then by Lemma E.2 and the
 755 inductive hypothesis, it samples from $q_{t_j}(\cdot \mid \hat{x}_{t_j})$. Alternatively, suppose the sampler performs some
 756 sequence of ℓ unmasking steps before performing an insertion step. Then by further inducting on ℓ ,
 757 we conclude by Lemma E.1 that the sampler eventually outputs a sample from $q_{t_j}(\cdot \mid \hat{x}_{t_j})$.

758 Finally, the theorem follows from the special case where $t_j = 0$ and $\hat{x}_{t_j} = \varepsilon$. □

759 **E.5 Proof of Lemma E.1**

760 *Proof.* Fix any index $i \in \text{mask}(x)$. The marginal $q_t^i(\cdot \mid x)$ is given by

$$\Pr_{(x_1, S) \mid x_t = x} [(x_1)_{s_i} = v] = \frac{\sum_{x_1, S: (x_1)_{s_i} = v} p(x) \cdot (1 - \alpha_t)^{\text{len}(x_1) - \text{len}(x)} \cdot \#\{S : x \subseteq (x_1) \mid_S\}}{\sum_{x_1} p(x_1) (1 - \alpha_t)^{\text{len}(x_1) - \text{len}(x)} \cdot \#\{S : x \subseteq (x_1) \mid_S\}}. \quad (10)$$

761 The posterior $q_t(\cdot | x[x^i \leftarrow v])$ is given by

$$q_t(x^* | x[x^i \leftarrow v]) = \frac{p(x^*) \cdot (1 - \alpha_t)^{\text{len}(x^*) - \text{len}(x)} \cdot \#\{S : x[x^i \leftarrow v] \subseteq x^*|_S\}}{\sum_{x_1} p(x_1) \cdot (1 - \alpha_t)^{\text{len}(x_1) - \text{len}(x)} \cdot \#\{S : x[x^i \leftarrow v] \subseteq (x_1)|_S\}}. \quad (11)$$

762 Note that the numerator of Eq. (10) and the denominator of Eq. (11) are the same. So conditioned on
763 unmasking index i , the above procedure outputs x^* with probability

$$\begin{aligned} & \sum_z \Pr_{(x_1, S) | x_t = x} [(x_1)|_{s_i} = v] \cdot q_t(x^* | x[x^i \leftarrow v]) \\ &= \frac{p(x^*) \cdot (1 - \alpha_t)^{\text{len}(x^*) - \text{len}(x)} \cdot \sum_z \#\{S : x[x^i \leftarrow v] \subseteq x^*|_S\}}{\sum_{x_1} p(x_1) (1 - \alpha_t)^{\text{len}(x_1) - \text{len}(x)} \cdot \#\{S : x \subseteq (x_1)|_S\}} \\ &= q_t(x^* | x). \end{aligned}$$

764 This holds conditioned on unmasking any $i \in \text{mask}(x)$, so regardless of the choice of π over such
765 positions, we will sample from the correct distribution $q_t(\cdot | x)$. \square

766 E.6 Proof of Lemma E.2

767 *Proof.* It suffices to show that the rate matrix satisfies the Kolmogorov *backward* equation

$$\partial_t q_t(x^* | x) = - \sum_{i=0}^{\text{len}(x)} R_t^{\text{ins}}(x, x \triangleleft_i \mathbf{m}) q_t(x^* | x \triangleleft_i \mathbf{m}) - R_t^{\text{ins}}(x, x) q_t(x^* | x).$$

768 First note that the rate $R_t^{\text{ins}}(x, x \triangleleft_i \mathbf{m})$ can be expressed as

$$\frac{\sum_{x_1} p(x_1) \cdot (1 - \alpha_t)^{\text{len}(x_1) - \text{len}(x)} \cdot \sum_{S: x \subseteq (x_1)|_S} (s_{i+1} - s_i - 1)}{\sum_{x_1} p(x_1) \cdot (1 - \alpha_t)^{\text{len}(x_1) - \text{len}(x)} \cdot \#\{S : x \subseteq (x_1)|_S\}} \cdot \frac{\dot{\alpha}_t}{1 - \alpha_t}.$$

769 Furthermore, $\sum_i (s_{i+1} - s_i - 1) = \text{len}(x_1) - \text{len}(x)$, so

$$\begin{aligned} & \sum_i R_t^{\text{ins}}(x, x \triangleleft_i \mathbf{m}) \\ &= \frac{\sum_{x_1} p(x_1) \cdot (1 - \alpha_t)^{\text{len}(x_1) - \text{len}(x)} \cdot \#\{S : x \subseteq (x_1)|_S\} \cdot (\text{len}(x_1) - \text{len}(x))}{\sum_{x_1} p(x_1) \cdot (1 - \alpha_t)^{\text{len}(x_1) - \text{len}(x)} \cdot \#\{S : x \subseteq (x_1)|_S\}} \cdot \frac{\dot{\alpha}_t}{1 - \alpha_t}. \quad (12) \end{aligned}$$

770 Let us compute $\partial_t q_t(x^* | x)$:

$$\begin{aligned} & - \frac{p(x^*) \cdot (1 - \alpha_t)^{\text{len}(x^*) - \text{len}(x)} \cdot \#\{S : x \subseteq x^*|_S\} \cdot (\text{len}(x^*) - \text{len}(x))}{\sum_{x_1} p(x_1) \cdot (1 - \alpha_t)^{\text{len}(x_1) - \text{len}(x)} \cdot \#\{S : x \subseteq (x_1)|_S\}} \cdot \frac{\dot{\alpha}_t}{1 - \alpha_t} \\ & + \left[\frac{\sum_{x_1} p(x_1) \cdot (1 - \alpha_t)^{\text{len}(x_1) - \text{len}(x)} \cdot \#\{S : x \subseteq (x_1)|_S\} \cdot (\text{len}(x_1) - \text{len}(x))}{\sum_{x_1} p(x_1) \cdot (1 - \alpha_t)^{\text{len}(x_1) - \text{len}(x)} \cdot \#\{S : x \subseteq (x_1)|_S\}} \cdot \frac{\dot{\alpha}_t}{1 - \alpha_t} \right. \\ & \quad \left. \times \frac{p(x^*) \cdot (1 - \alpha_t)^{\text{len}(x^*) - \text{len}(x)} \cdot \#\{S : x \subseteq x^*|_S\}}{\sum_{x_1} p(x_1) \cdot (1 - \alpha_t)^{\text{len}(x_1) - \text{len}(x)} \cdot \#\{S : x \subseteq (x_1)|_S\}} \right] \quad (13) \end{aligned}$$

771 Note that by Eq. (12), the term in the parentheses in Eq. (13) is exactly

$$\sum_{i=0}^{\text{len}(x)} R_t^{\text{ins}}(x, x \triangleleft_i \mathbf{m}) q_t(x^* | x) = -R_t^{\text{ins}}(x, x) q_t(x^* | x),$$

772 It remains to verify that the first term in Eq. (13) is equal to $-\sum_i R_t^{\text{ins}}(x, x \triangleleft_i \mathbf{m}) q_t(x^* | x \triangleleft_i \mathbf{m})$.
773 To that end, we must show that

$$\begin{aligned} & \sum_{i=0}^{\text{len}(x)} \frac{\sum_{x_1} p(x_1) (1 - \alpha_t)^{\text{len}(x_1) - \text{len}(x)} \cdot \sum_{S: x \subseteq (x_1)|_S} (s_{i+1} - s_i - 1)}{\sum_{x_1} p(x_1) (1 - \alpha_t)^{\text{len}(x_1) - \text{len}(x)} \cdot \#\{S : x_t \triangleleft_i \mathbf{m} \subseteq (x_1)|_S\}} \cdot \#\{S : x \triangleleft_i \mathbf{m} \subseteq x^*|_S\} \\ &= \#\{S : x \subseteq (x^*)|_S\} \cdot (\text{len}(x^*) - \text{len}(x)) \quad (14) \end{aligned}$$

774 The key combinatorial step is as follows. For any x_1 in the support of p , consider a subset S for which
 775 $x \subseteq (x_1)|_S$. Note that for every such S , we can uniquely associate exactly $s_{i+1} - s_i - 1$ different
 776 subsets S' of size $|S| + 1$ for which $x \triangleleft_i \mathbf{m} \subseteq (x_1)|_{S'}$. Therefore, $\sum_{S: x \subseteq (x_1)|_S} (s_{i+1} - s_i - 1) =$
 777 $\#\{S : x \triangleleft_i \mathbf{m} \subseteq (x_1)|_S\}$, and the left-hand side of Eq. (14) thus becomes

$$\sum_{i=0}^{\text{len}(x)} \#\{S : x \triangleleft_i \mathbf{m} \subseteq x^*|_S\} = \sum_{i=0}^{\text{len}(x)} \sum_{S: x \subseteq x^*|_S} (s_{i+1} - s_i - 1) = \sum_{S: x \subseteq (x_1)|_S} (\text{len}(x^*) - \text{len}(x_t)),$$

778 which completes the proof of Eq. (14). \square

779 F Experimental details

780 F.1 Pretraining on OpenWebText

781 **Dataset preparation.** As mentioned in Section 5.1, to obtain a variable-length dataset, we split
 782 OpenWebText articles paragraph-wise using the GPT-2 tokenizer Radford and Wu [2019]. This can
 783 be implemented by locating the token index corresponding to the delimiter `\.`. Sequences longer
 784 than 1024 tokens are then chunked, yielding a variable-length dataset with maximum sequence length
 785 1024.

786 **FlexAttention.** To handle variable-length sequences during training, we pad each batch to the
 787 maximum sequence length. In MDM, by design, pad tokens also enter the model input. In contrast,
 788 FlexMDM is designed to receive only clean or mask tokens as inputs. Ideally, QKV attention should
 789 not attend to pad tokens; however, current FlashAttention [Dao et al., 2022] does not support this for
 790 *non-causal* attention (our setup of interest). We therefore adapt FlexAttention [Dong et al., 2024].
 791 A side benefit is improved training speed, since FlexMDM performs attention over fewer tokens
 792 than MDM’s full-sequence attention. Note that in the LLaDA experiment, we did not implement
 793 this optimization; pad tokens can therefore attend to other tokens, though we expect the impact to be
 794 negligible.

795 **Training configuration.** As in Section 5.1, we model FlexMDM with a DiT Peebles and Xie
 796 [2023] backbone and embed the insertion schedule α_t . For MDM, we use the same DiT configuration
 797 with time step embedding but without the softplus scalar head. Transformer configuration is: `hidden`
 798 `size:768`, `heads:12`, `conditional dimension:128`, `dropout:0.05`, `layers:12`. We train
 799 both models on 16 H100 GPUs with a global batch size of 1024 and max training iteration $1M$. We
 800 use the AdamW [Loshchilov and Hutter, 2017] optimizer with weight decay 0.03, learning rate $3e-4$,
 801 2000 warmup steps, and an EMA factor of 0.9999. Additionally, we use low-discrepancy time-step
 802 sampling to reduce variance: one t is drawn uniformly from each interval $[i/T, (i+1)/T]_{i=0}^{T-1}$, as in
 803 prior MDM training [Shi et al., 2024].

804 **Metric.** For evaluation, we take sequences generated by both models and retain the clean tokens by
 805 removing padding (e.g., the leading pad token). We adopt LLaMA-2-7B [Touvron et al., 2023] as the
 806 reference model to compute likelihoods.

807 F.2 Pretraining on the Maze planning task

808 **Task construction.** We generate mazes with a fully random, recursive division procedure (Listing 1),
 809 on a 41×41 grid where some cells are invalid. As described in Section 5.1, we use a subgoal-
 810 conditioned planning task: the model is given a sequence of subgoals (g_1, \dots, g_K) and must produce
 811 a valid path that connects them in order. To construct training examples for a given K , we sample
 812 15000 start–goal pairs, compute the shortest path for each pair via breadth-first search, and then add
 813 controlled detours to obtain up to 10 distinct valid paths per pair. Subgoals are formed by selecting
 814 $K - 2$ intermediate cells uniformly at random along a chosen path (start and goal are already fixed).
 815 We use 95% of the pairs for training and hold out 5% for validation to evaluate generalization to
 816 unseen pairs and subgoal sets.

817 **Training data construction.** For MDM, the training sequence is $((g_1, \dots, g_K) \text{ [SEP] Path})$,
 818 where [SEP] is a special separator and Path denotes the tokenized path. During training, the

819 prompt (g_1, \dots, g_K) bypasses the interpolant so that, at inference time, the model can condition
 820 on (g_1, \dots, g_K) [SEP] and generate the path. For FlexMDM, we use an interpolant in which the
 821 subgoal tokens are exempt from the process in (6); that is, tokens corresponding to each g_i are kept
 822 clean at all times. This also enables generation to start from (g_1, \dots, g_K) .

823 **Training configuration.** We use the same design as in the OpenWebText pretraining but with
 824 a smaller model: hidden size:256, heads:8, conditional dimension:128, dropout:0.1,
 825 layers:8. Both models have approximately 8.5M parameters. We train on 4 A100 GPUs with a
 826 global batch size of 256 for up to 100 epochs. We use AdamW [Loshchilov and Hutter, 2017] with
 827 weight decay 0.01, learning rate 1×10^{-4} , 20 warmup epochs, and an EMA factor of 0.9999.

828 **Metric.** Given the final conditionally generated sequence, we report success rate under two criteria:
 829 (1) all visited cells are valid (no collisions with blocked cells), and (2) the path connects the subgoals
 830 consecutively in order.

Listing 1: Code for the maze Construction

```

831 #
832 -----
833
834 # RECURSIVE DIVISION (perfect maze)
835 -----
836 #
837 -----
838
839 def _divide(g, top, left, h, w):
840     if h <= 2 or w <= 2:
841         return
842     horizontal = w < h # split the longer dimension
843     if horizontal:
844         y = random.randrange(top + 1, top + h - 1, 2)
845         gap = random.randrange(left, left + w, 2)
846         g[y, left:left + w] = 1
847         g[y, gap] = 0
848         _divide(g, top, left, y - top, w)
849         _divide(g, y + 1, left, top + h - y - 1, w)
850     else:
851         x = random.randrange(left + 1, left + w - 1, 2)
852         gap = random.randrange(top, top + h, 2)
853         g[top:top + h, x] = 1
854         g[gap, x] = 0
855         _divide(g, top, left, h, x - left)
856         _divide(g, top, x + 1, h, left + w - x - 1)
857
858 #
859 -----
860
861 # WRAPPER with extra doorways -----
862 #
863 -----
864
865 def division_maze(m, n, seed=2025, extra_door_frac=0.5):
866     """
867     m, n # size in *cells* (not bitmap coords)
868     seed # int or None
869     extra_door_frac # 0, perfect maze; >0 flicks more doors (loops)
870     """
871     random.seed(seed)
872     H, W = 2 * m + 1, 2 * n + 1
873     g = np.zeros((H, W), dtype=np.uint8)
874     g[0, :], g[H - 1, :], g[:, 0], g[:, W - 1] = 1, 1, 1, 1
875
876     _divide(g, 1, 1, H - 2, W - 2)

```

```

878
879 # ----- optional imperfection: punch extra doorways
880 -----
881 if extra_door_frac > 0:
882     candidates = []
883     for r in range(1, H - 1):
884         for c in range(1, W - 1):
885             if g[r, c] == 1:
886                 # Check if wall separates two passages
887                 # orthogonally
888                 if g[r - 1, c] == 0 and g[r + 1, c] == 0:
889                     candidates.append((r, c))
890                 elif g[r, c - 1] == 0 and g[r, c + 1] == 0:
891                     candidates.append((r, c))
892     k = int(len(candidates) * extra_door_frac)
893     for (r, c) in random.sample(candidates, k=k):
894         g[r, c] = 0
895     return g

```

896 F.3 Weight Initialization training from LLaDA

897 In this appendix, we describe the procedure for adapting the pretrained LLaDa 8B base model into
898 the FlexMDM, a natural choice given the design FlexMDM builds incrementally upon the design of
899 MDM.

900 **Training configuration** The LLaDa model uses a bidirectional transformer architecture, relying
901 on the fact that masked diffusion model does not necessitate a time embedding Zheng et al. [2024].
902 We reinject a time-embedding via AdaLN Peebles and Xie [2023] as it is necessary for maintaining
903 insertion expectation.

904 For every attention layer, q_proj , k_proj , v_proj are then finetuned through LoRA Hu et al. [2022].
905 The unmasking posterior layer (a linear layer) is also finetuned through LoRA, as it is used for a
906 semantically different task in FlexMDM where the unmasking has to account for shifting of tokens
907 due to insertion. We highlight that it is crucial to finetune this layer to obtain readable generations.
908 For both cases, LoRA is used with hyperparameters $r = 128$, $\alpha = 128$ and dropout of rate 0.1.

909 Training is done for 200000 gradient steps with batch size 64 on 16 H100s which takes approximately
910 three days. Optimisation is done through the AdamW Loshchilov and Hutter [2017] with learning
911 rate $1e - 4$ and weight decay of 0.1 with a cosine warm restart scheduler.

912 For instruction fine-tuning, in the case of GSM8k, both the base model and FlexMDM are instruction
913 finetuned for 1000 epochs, whereas in the case of infilling, both models are finetuned for 50 epochs.
914 Training configurations are the same as pretraining.

915 **Sampling.** For FlexMDM, we adopt the Top-K probability confidence-based sampling with sliding
916 window with $\gamma_1 = 5.0$, $\gamma_2 = 64$, for details of the algorithm see Appendix E. For MDM, we
917 report the best of a semiautoregressive confidence-based sampling strategy and confidence-based
918 strategy following Nie et al. [2024]. In the case for FlexMDM, We highlight that using a sliding
919 window significantly improves performance, which is surprisingly not the case for MDM with
920 semi-autoregressive sampling.

921 **NeurIPS Paper Checklist**

922 The checklist is designed to encourage best practices for responsible machine learning research,
923 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
924 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should
925 follow the references and follow the (optional) supplemental material. The checklist does NOT count
926 towards the page limit.

927 Please read the checklist guidelines carefully for information on how to answer these questions. For
928 each question in the checklist:

- 929 • You should answer [Yes], [No], or [NA].
- 930 • [NA] means either that the question is Not Applicable for that particular paper or the
931 relevant information is Not Available.
- 932 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

933 **The checklist answers are an integral part of your paper submission.** They are visible to the
934 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it
935 (after eventual revisions) with the final version of your paper, and its final version will be published
936 with the paper.

937 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
938 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a
939 proper justification is given (e.g., "error bars are not reported because it would be too computationally
940 expensive" or "we were unable to find the license for the dataset we used"). In general, answering
941 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we
942 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
943 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
944 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification
945 please point to the section(s) where related material for the question can be found.

946 IMPORTANT, please:

- 947 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”,**
- 948 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 949 • **Do not modify the questions and only use the provided macros for your answers.**

950 **1. Claims**

951 Question: Do the main claims made in the abstract and introduction accurately reflect the
952 paper’s contributions and scope?

953 Answer: **[TODO]**

954 Justification: **[TODO]**

955 Guidelines:

- 956 • The answer NA means that the abstract and introduction do not include the claims
957 made in the paper.
- 958 • The abstract and/or introduction should clearly state the claims made, including the
959 contributions made in the paper and important assumptions and limitations. A No or
960 NA answer to this question will not be perceived well by the reviewers.
- 961 • The claims made should match theoretical and experimental results, and reflect how
962 much the results can be expected to generalize to other settings.
- 963 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
964 are not attained by the paper.

965 **2. Limitations**

966 Question: Does the paper discuss the limitations of the work performed by the authors?

967 Answer: **[TODO]**

968 Justification: **[TODO]**

969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [TODO]

Justification: [TODO]

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [TODO]

Justification: [TODO]

Guidelines:

- The answer NA means that the paper does not include experiments.

- 1020 • If the paper includes experiments, a No answer to this question will not be perceived
1021 well by the reviewers: Making the paper reproducible is important, regardless of
1022 whether the code and data are provided or not.
- 1023 • If the contribution is a dataset and/or model, the authors should describe the steps taken
1024 to make their results reproducible or verifiable.
- 1025 • Depending on the contribution, reproducibility can be accomplished in various ways.
1026 For example, if the contribution is a novel architecture, describing the architecture fully
1027 might suffice, or if the contribution is a specific model and empirical evaluation, it may
1028 be necessary to either make it possible for others to replicate the model with the same
1029 dataset, or provide access to the model. In general, releasing code and data is often
1030 one good way to accomplish this, but reproducibility can also be provided via detailed
1031 instructions for how to replicate the results, access to a hosted model (e.g., in the case
1032 of a large language model), releasing of a model checkpoint, or other means that are
1033 appropriate to the research performed.
- 1034 • While NeurIPS does not require releasing code, the conference does require all submis-
1035 sions to provide some reasonable avenue for reproducibility, which may depend on the
1036 nature of the contribution. For example
 - 1037 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
1038 to reproduce that algorithm.
 - 1039 (b) If the contribution is primarily a new model architecture, the paper should describe
1040 the architecture clearly and fully.
 - 1041 (c) If the contribution is a new model (e.g., a large language model), then there should
1042 either be a way to access this model for reproducing the results or a way to reproduce
1043 the model (e.g., with an open-source dataset or instructions for how to construct
1044 the dataset).
 - 1045 (d) We recognize that reproducibility may be tricky in some cases, in which case
1046 authors are welcome to describe the particular way they provide for reproducibility.
1047 In the case of closed-source models, it may be that access to the model is limited in
1048 some way (e.g., to registered users), but it should be possible for other researchers
1049 to have some path to reproducing or verifying the results.

1050 5. Open access to data and code

1051 Question: Does the paper provide open access to the data and code, with sufficient instruc-
1052 tions to faithfully reproduce the main experimental results, as described in supplemental
1053 material?

1054 Answer: **[TODO]**

1055 Justification: **[TODO]**

1056 Guidelines:

- 1057 • The answer NA means that paper does not include experiments requiring code.
- 1058 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
1059 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1060 • While we encourage the release of code and data, we understand that this might not be
1061 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
1062 including code, unless this is central to the contribution (e.g., for a new open-source
1063 benchmark).
- 1064 • The instructions should contain the exact command and environment needed to run to
1065 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
1066 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1067 • The authors should provide instructions on data access and preparation, including how
1068 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1069 • The authors should provide scripts to reproduce all experimental results for the new
1070 proposed method and baselines. If only a subset of experiments are reproducible, they
1071 should state which ones are omitted from the script and why.
- 1072 • At submission time, to preserve anonymity, the authors should release anonymized
1073 versions (if applicable).

- 1074 • Providing as much information as possible in supplemental material (appended to the
1075 paper) is recommended, but including URLs to data and code is permitted.

1076 **6. Experimental setting/details**

1077 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
1078 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
1079 results?

1080 Answer: **[TODO]**

1081 Justification: **[TODO]**

1082 Guidelines:

- 1083 • The answer NA means that the paper does not include experiments.
1084 • The experimental setting should be presented in the core of the paper to a level of detail
1085 that is necessary to appreciate the results and make sense of them.
1086 • The full details can be provided either with the code, in appendix, or as supplemental
1087 material.

1088 **7. Experiment statistical significance**

1089 Question: Does the paper report error bars suitably and correctly defined or other appropriate
1090 information about the statistical significance of the experiments?

1091 Answer: **[TODO]**

1092 Justification: **[TODO]**

1093 Guidelines:

- 1094 • The answer NA means that the paper does not include experiments.
1095 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
1096 dence intervals, or statistical significance tests, at least for the experiments that support
1097 the main claims of the paper.
1098 • The factors of variability that the error bars are capturing should be clearly stated (for
1099 example, train/test split, initialization, random drawing of some parameter, or overall
1100 run with given experimental conditions).
1101 • The method for calculating the error bars should be explained (closed form formula,
1102 call to a library function, bootstrap, etc.)
1103 • The assumptions made should be given (e.g., Normally distributed errors).
1104 • It should be clear whether the error bar is the standard deviation or the standard error
1105 of the mean.
1106 • It is OK to report 1-sigma error bars, but one should state it. The authors should
1107 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
1108 of Normality of errors is not verified.
1109 • For asymmetric distributions, the authors should be careful not to show in tables or
1110 figures symmetric error bars that would yield results that are out of range (e.g. negative
1111 error rates).
1112 • If error bars are reported in tables or plots, The authors should explain in the text how
1113 they were calculated and reference the corresponding figures or tables in the text.

1114 **8. Experiments compute resources**

1115 Question: For each experiment, does the paper provide sufficient information on the com-
1116 puter resources (type of compute workers, memory, time of execution) needed to reproduce
1117 the experiments?

1118 Answer: **[TODO]**

1119 Justification: **[TODO]**

1120 Guidelines:

- 1121 • The answer NA means that the paper does not include experiments.
1122 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
1123 or cloud provider, including relevant memory and storage.

- 1124 • The paper should provide the amount of compute required for each of the individual
1125 experimental runs as well as estimate the total compute.
1126 • The paper should disclose whether the full research project required more compute
1127 than the experiments reported in the paper (e.g., preliminary or failed experiments that
1128 didn't make it into the paper).

1129 9. Code of ethics

1130 Question: Does the research conducted in the paper conform, in every respect, with the
1131 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1132 Answer: **[TODO]**

1133 Justification: **[TODO]**

1134 Guidelines:

- 1135 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 1136 • If the authors answer No, they should explain the special circumstances that require a
1137 deviation from the Code of Ethics.
- 1138 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
1139 eration due to laws or regulations in their jurisdiction).

1140 10. Broader impacts

1141 Question: Does the paper discuss both potential positive societal impacts and negative
1142 societal impacts of the work performed?

1143 Answer: **[TODO]**

1144 Justification: **[TODO]**

1145 Guidelines:

- 1146 • The answer NA means that there is no societal impact of the work performed.
- 1147 • If the authors answer NA or No, they should explain why their work has no societal
1148 impact or why the paper does not address societal impact.
- 1149 • Examples of negative societal impacts include potential malicious or unintended uses
1150 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
1151 (e.g., deployment of technologies that could make decisions that unfairly impact specific
1152 groups), privacy considerations, and security considerations.
- 1153 • The conference expects that many papers will be foundational research and not tied
1154 to particular applications, let alone deployments. However, if there is a direct path to
1155 any negative applications, the authors should point it out. For example, it is legitimate
1156 to point out that an improvement in the quality of generative models could be used to
1157 generate deepfakes for disinformation. On the other hand, it is not needed to point out
1158 that a generic algorithm for optimizing neural networks could enable people to train
1159 models that generate Deepfakes faster.
- 1160 • The authors should consider possible harms that could arise when the technology is
1161 being used as intended and functioning correctly, harms that could arise when the
1162 technology is being used as intended but gives incorrect results, and harms following
1163 from (intentional or unintentional) misuse of the technology.
- 1164 • If there are negative societal impacts, the authors could also discuss possible mitigation
1165 strategies (e.g., gated release of models, providing defenses in addition to attacks,
1166 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
1167 feedback over time, improving the efficiency and accessibility of ML).

1168 11. Safeguards

1169 Question: Does the paper describe safeguards that have been put in place for responsible
1170 release of data or models that have a high risk for misuse (e.g., pretrained language models,
1171 image generators, or scraped datasets)?

1172 Answer: **[TODO]**

1173 Justification: **[TODO]**

1174 Guidelines:

- 1175 • The answer NA means that the paper poses no such risks.

- 1176 • Released models that have a high risk for misuse or dual-use should be released with
1177 necessary safeguards to allow for controlled use of the model, for example by requiring
1178 that users adhere to usage guidelines or restrictions to access the model or implementing
1179 safety filters.
- 1180 • Datasets that have been scraped from the Internet could pose safety risks. The authors
1181 should describe how they avoided releasing unsafe images.
- 1182 • We recognize that providing effective safeguards is challenging, and many papers do
1183 not require this, but we encourage authors to take this into account and make a best
1184 faith effort.

1185 12. Licenses for existing assets

1186 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1187 the paper, properly credited and are the license and terms of use explicitly mentioned and
1188 properly respected?

1189 Answer: **[TODO]**

1190 Justification: **[TODO]**

1191 Guidelines:

- 1192 • The answer NA means that the paper does not use existing assets.
- 1193 • The authors should cite the original paper that produced the code package or dataset.
- 1194 • The authors should state which version of the asset is used and, if possible, include a
1195 URL.
- 1196 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1197 • For scraped data from a particular source (e.g., website), the copyright and terms of
1198 service of that source should be provided.
- 1199 • If assets are released, the license, copyright information, and terms of use in the
1200 package should be provided. For popular datasets, `paperswithcode.com/datasets`
1201 has curated licenses for some datasets. Their licensing guide can help determine the
1202 license of a dataset.
- 1203 • For existing datasets that are re-packaged, both the original license and the license of
1204 the derived asset (if it has changed) should be provided.
- 1205 • If this information is not available online, the authors are encouraged to reach out to
1206 the asset's creators.

1207 13. New assets

1208 Question: Are new assets introduced in the paper well documented and is the documentation
1209 provided alongside the assets?

1210 Answer: **[TODO]**

1211 Justification: **[TODO]**

1212 Guidelines:

- 1213 • The answer NA means that the paper does not release new assets.
- 1214 • Researchers should communicate the details of the dataset/code/model as part of their
1215 submissions via structured templates. This includes details about training, license,
1216 limitations, etc.
- 1217 • The paper should discuss whether and how consent was obtained from people whose
1218 asset is used.
- 1219 • At submission time, remember to anonymize your assets (if applicable). You can either
1220 create an anonymized URL or include an anonymized zip file.

1221 14. Crowdsourcing and research with human subjects

1222 Question: For crowdsourcing experiments and research with human subjects, does the paper
1223 include the full text of instructions given to participants and screenshots, if applicable, as
1224 well as details about compensation (if any)?

1225 Answer: **[TODO]**

1226 Justification: **[TODO]**

1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.