

# BEYOND MASKS: EFFICIENT, FLEXIBLE DIFFUSION LANGUAGE MODELS VIA DELETION-INSERTION PROCESSES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

While Masked Diffusion Language Models (MDLMs) relying on token masking and unmasking have shown promise in language modeling, their computational efficiency and generation flexibility remain constrained by the masking paradigm. In this paper, we propose Deletion-Insertion Diffusion language models (DID) that rigorously formulate token deletion and insertion as discrete diffusion processes, replacing the masking and unmasking processes in current MDLMs. DID improves training and inference efficiency by eliminating two major sources of computational overhead in MDLMs: the computations on non-informative 1) `<MASK>` tokens inherent to the paradigm, and 2) `<PAD>` tokens introduced in variable-length settings. Furthermore, DID offers greater flexibility by: 1) natively supporting variable-length sequences without requiring fixed-length padding, and 2) an intrinsic self-correction mechanism during generation due to insertion that dynamically adjusts token positions. To train DID, we design a score-based approach that assigns scores to token insertion operations and derive appropriate training objectives. The objectives involve subsequence counting problems, which we efficiently solve via a parallelized dynamic programming algorithm. Our experiments across fixed and variable-length settings demonstrate the advantage of DID over baselines of MDLMs and existing insertion-based LMs, in terms of modeling performance, sampling quality, and training/inference speed.

## 1 INTRODUCTION

Diffusion language models (DLMs) (Austin et al., 2021; Campbell et al., 2022; Lou et al., 2024) have rapidly emerged as a powerful paradigm for language modeling, offering a compelling alternative to the dominant autoregressive (AR) approach. They offer distinct advantages, including bidirectional context modeling and the potential for parallel decoding. Within this domain, a body of work on Masked Diffusion Language Models (MDLMs) (Nie et al., 2024; 2025; Ou et al., 2025; Sahoo et al., 2024; Shi et al., 2024) is the most widely studied. These models operate through a forward process that progressively corrupts each token into an absorbing state `<MASK>` and a backward process that reconstructs the original sequence by iteratively unmasking tokens from a fully masked sequence, with a fixed sequence length during the diffusion process.

Despite their success, MDLMs are fundamentally limited by their fixed sequence length. The first issue lies in their restricted generation flexibility, which leads to challenges in modeling variable lengths and performing self-correction: once a token is unmasked, its content and position become fixed, thereby risking error accumulation in a similar sense to autoregressive models. The second issue is their substantial computational inefficiency, as the model must repeatedly process full-length sequences. Under the typical log-linear noise schedule, about half of the FLOPs are allocated to the non-informative `<MASK>` tokens during both training and inference. Further, if MDLMs are applied to variable-length sequences, their fixed-length nature demands padding to the same length (Nie et al., 2024; 2025; Wu et al., 2025b; Gong et al., 2024) (Fig. 1a), allocating extra FLOPs to the non-informative `<PAD>` tokens. This means generating a shorter sequence is not faster.

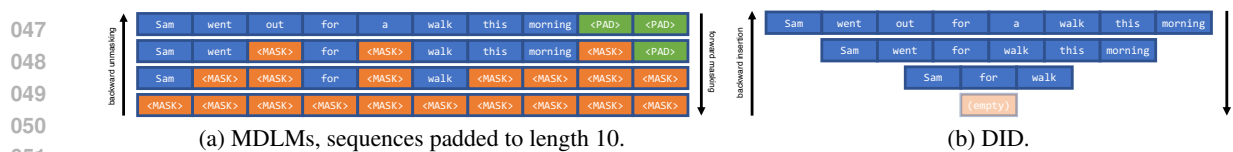


Figure 1: Conceptual diagram of MDLMs compared to Deletion-Insertion Diffusion language models (DID).

To address these issues, we propose Deletion-Insertion Diffusion language models (DID), a novel discrete diffusion paradigm that fundamentally differs from MDLMs. DID replaces the masking-unmasking processes in MDLMs with deletion-insertion processes (Fig. 1b). Specifically, tokens are progressively deleted in the forward process until the sequence is empty; in the backward process, generation starts from an empty sequence and iteratively inserts tokens until a complete sequence is reconstructed. DID eliminates the <MASK> and <PAD> tokens used in MDLMs, saving FLOPs and improving computational efficiency. Regarding generation flexibility, DID natively supports variable-length data, and as an insertion-based language model, features an intrinsic self-correction mechanism that dynamically adjusts token positions during generation.

We implement DID by addressing several non-trivial design and training challenges. First, we rigorously formulate the deletion and insertion processes within the discrete diffusion framework, and develop a score-based approach built upon the Denoising Score Entropy (DSE) (Lou et al., 2024) objective to train DID. Concretely, we define an insertion score that models the probability of inserting any token at any position of a sequence at a given time interval, and derive a corresponding Denoising Insertion Score Entropy (DISE) training objective. The DISE objective is based on a ratio of subsequence counts in the clean data after and before an insertion, which serves as the training target for the insertion score. To efficiently compute the ratio, we develop a parallelized dynamic programming algorithm that exploits parallelism. Moreover, we demonstrate that under the fixed-length setting of MDLMs, the DISE objective can be further simplified to a form resembling cross-entropy, further improving parameterization and learning of the insertion score.

Comprehensive experiments demonstrate the effectiveness of DID in enhancing efficiency and flexibility. In fixed-length language modeling benchmarks, DID achieves superior performance compared to strong MDLM baselines (e.g., RADD (Ou et al., 2025)) when aligned by computational budget (FLOPs) (Tab. 1). Compared to MDLM baselines, DID could accelerate training by up to  $1.82\times$  and  $3.10\times$  (Tab. 3, 5) and inference by up to  $1.58\times$  and  $3.79\times$  (Tab. 2, 4), for models trained on fixed-length and variable-length datasets, respectively. Moreover, DID shows strong performance in variable-length settings, outperforming MDLMs and existing insertion-based LMs in sampling quality and consistency with data length distribution (Tab. 4, Fig. 2).

In summary, our main contributions are as follows:

- We propose DID, a novel diffusion LM based on deletion-insertion processes that eliminates the use of <MASK> and <PAD> tokens, improving the computational efficiency and generation flexibility of DLMs.
- We develop DISE, a score-based training objective, and an efficient parallel dynamic programming implementation that enables effective learning of DID’s insertion process for language generation.
- Our experiments demonstrate the superior efficiency and flexibility of DID over baselines of MDLMs and other insertion-based LMs on language/length modeling, generation quality, and training/inference speed.

## 2 PRELIMINARIES AND RELATED WORKS

### 2.1 CONTINUOUS-TIME DISCRETE DIFFUSION

A continuous-time discrete diffusion model consists of a forward noising process and a backward denoising process, both continuous-time Markov chains. In the forward process, the training samples are progressively corrupted into pure noise. The model aims to learn its backward process that inverts this corruption, and generate new samples by sampling through the backward process from noise.

**Continuous-time Markov chain.** We consider a discrete state space  $\mathcal{X}$ . A continuous-time Markov chain (CTMC) is a process  $\mathbf{x}_t$  on  $\mathcal{X}$  with  $t \in [0, 1]$ , starting from an initial data distribution  $p_0(\mathbf{x}_0)$ . A CTMC is characterized by a time-dependent transition rate matrix  $Q_t \in \mathbb{R}^{|\mathcal{X}| \times |\mathcal{X}|}$ . For distinct states  $\mathbf{x}_t, \mathbf{y} \in \mathcal{X}$ ,  $Q_t(\mathbf{x}_t, \mathbf{y}) \geq 0$  defines the instantaneous transition rate from  $\mathbf{x}_t$  to  $\mathbf{y}$ . This means, at an infinitesimal time interval  $[t, t + \Delta t]$ , the transition probability is given by:

$$p_{t+\Delta t|t}(\mathbf{y}|\mathbf{x}_t) = \delta(\mathbf{x}_t, \mathbf{y}) + Q_t(\mathbf{x}_t, \mathbf{y})\Delta t, \quad (1)$$

where  $\delta$  is the Kronecker delta. In other words, the evolution of the marginal distribution  $p_t \in \Delta_{|\mathcal{X}|}$  follows the Kolmogorov forward equation  $\frac{dp_t}{dt} = p_t Q_t$ . Note that the diagonal entries of  $Q_t$  should satisfy  $Q_t(\mathbf{x}_t, \mathbf{x}_t) = -\sum_{\mathbf{x}_t \neq \mathbf{y}} Q_t(\mathbf{x}_t, \mathbf{y})$  to ensure the weights of  $p_t$  add up to 1.

**Determining the forward process.** In the forward process, a common parameterization (Campbell et al., 2022) of the transition rate matrix is  $Q_t = \sigma(t)Q$ , where  $\sigma(t)$  is a scalar noise schedule and  $Q$  is a constant rate matrix determined by the model design. Taking this to the Kolmogorov forward equation, one can analytically solve the marginal distributions by  $p_t = p_s P_{t|s}$ , where  $P_{t|s} = \exp((\bar{\sigma}(t) - \bar{\sigma}(s))Q)$  is the transition probability matrix from time  $s$  to time  $t$ , and  $\bar{\sigma}(t) = \int_0^t \sigma(\tau) d\tau$ .

**Learning the backward process.** It is known that the time reversal of this process is also a CTMC, with its infinitesimal transition probability similar to Eq. 1:

$$p_{t-\Delta t|t}(\mathbf{y}|\mathbf{x}_t) = \delta(\mathbf{x}_t, \mathbf{y}) + \tilde{Q}_t(\mathbf{x}_t, \mathbf{y})\Delta t. \quad (2)$$

The reverse transition rate matrix  $\tilde{Q}_t$  is associated to its forward counterpart  $Q_t$  by the identity  $\tilde{Q}_t(\mathbf{x}_t, \mathbf{y}) = Q_t(\mathbf{y}, \mathbf{x}_t) s(\mathbf{x}_t, t)_{\mathbf{y}}$  for  $\mathbf{x}_t \neq \mathbf{y}$ , and  $\tilde{Q}_t(\mathbf{x}_t, \mathbf{x}_t) = -\sum_{\mathbf{x}_t \neq \mathbf{y}} \tilde{Q}_t(\mathbf{x}_t, \mathbf{y})$  (Lou et al., 2024). Here,  $s(\mathbf{x}_t, t)_{\mathbf{y}} = p_t(\mathbf{y})/p_t(\mathbf{x}_t)$  is the **concrete score**, which is generally intractable and commonly approximated by a parameterized network  $s_\theta(\mathbf{x}_t, t)_{\mathbf{y}}$  trained using the Denoising Score Entropy (DSE) objective (Lou et al., 2024), a negative evidence lower bound (ELBO) for discrete diffusion models (details in Appendix D.1):

$$\mathcal{L}_\theta^{\text{DSE}}(\mathbf{x}_0) = \mathbb{E}_{t, \mathbf{x}_t} \sum_{\mathbf{y} \neq \mathbf{x}_t} Q_t(\mathbf{y}, \mathbf{x}_t) \left[ s_\theta(\mathbf{x}_t, t)_{\mathbf{y}} - \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \log s_\theta(\mathbf{x}_t, t)_{\mathbf{y}} + K \left( \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \right) \right], \quad (3)$$

where the expectation is taken over  $t \sim \text{Unif}([0, 1])$  and  $\mathbf{x}_t \sim p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)$ , and  $K(a) = a(\log a - 1)$ .

## 2.2 INSERTION-BASED LANGUAGE MODELS

Classical insertion-based models, such as the Insertion Transformer (Stern et al., 2019), Levenshtein Transformer (Gu et al., 2019), and InsNet (Lu et al., 2022), pioneered sequence generation via iterative insertion or edit operations. These models demonstrate the potential of flexible decoding orders and parallel decoding. However, their training objectives are defined at the level of local edit policies, rather than arising from a probabilistic modeling of the global data distribution.

Recently, insertion operations have been integrated into discrete diffusion and flow matching. Flexible Masked Diffusion Models (FlexMDMs) (Kim et al., 2025) augment MDMs with a learned insertion expectation to insert additional `<MASK>` tokens during generation; this enables variable-length generation but still relies on the masking-unmasking paradigm. Edit Flows (Havasi et al., 2025) instead define a CTMC directly over variable-length sequences via insertion, deletion, and substitution edits. To make sequence-level flow matching tractable, Edit Flows augment the state space with auxiliary edit-path variables and estimate the training objective by Monte Carlo sampling of these paths. While this preserves theoretical correctness, it adds implementation complexity and an additional source of stochasticity beyond the randomness of the CTMC forward process. Insertion Language Models (ILMs) (Patel et al., 2025) draw inspiration from diffusion approaches and formulate a CTMC forward process that deletes tokens to learn a backward insertion process.

141 However, ILMs are not diffusion models and cannot learn the true backward insertion process since their  
 142 training objective is more heuristic rather than likelihood-bounded. They also suffer from several practical  
 143 limitations; for instance, they can only insert one token per step and require another network to determine  
 144 when to stop the generation.

145 DID bridges the gap between flexible insertion-based generation and rigorous discrete diffusion. Unlike  
 146 classical insertion models and ILMs, DID is well-grounded in a continuous-time diffusion framework:  
 147 deletion and insertion form the forward and backward CTMCs over the full variable-length sequence space,  
 148 and the DISE objective is inherited from the DSE objective to enable likelihood-bounded training. Unlike  
 149 Edit Flows, however, DID does not introduce auxiliary edit-path variables. Thanks to the independence  
 150 structure of the deletion process, the forward transition probabilities admit a closed-form expression in terms  
 151 of subsequence counts, and the subsequence count ratios that appear in DISE can be computed **exactly** via  
 152 parallel dynamic programming. Therefore, DID avoids the additional variance and engineering overhead  
 153 associated with sampling edit paths, while maintaining a principled, likelihood-bounded training objective  
 154 and completely eliminating the use of <MASK> and <PAD> tokens.

### 155 3 DID: DELETION-INSERTION DIFFUSION LANGUAGE MODELS

156 We propose DID to improve the efficiency and flexibility of diffusion language models. Instead of masking  
 157 and unmasking in MDLMs, DID reconstructs the diffusion processes with deletion and insertion. In this  
 158 section, we rigorously formulate the forward deletion process in Sec. 3.1, backward insertion process and  
 159 sampling algorithm in Sec. 3.2, develop a score-based approach to train DID in Sec. 3.3, discuss efficient  
 160 implementation supporting parallelism for DID training objectives in Sec. 3.4, and analyze the additional  
 161 optimization for the fixed-length data setting considered by MDLMs in Sec. 3.5.

#### 162 3.1 FORWARD PROCESS: DELETION

163 The forward process of DID is a CTMC on the state space  $\cup_{d=0}^{\infty} \mathcal{V}^d$  that gradually shortens the sequence  
 164 length by deleting tokens, thus equipping the model with variable-length ability. Similar to MDLMs, we  
 165 define this forward process through independent token-level deletions with rate  $\sigma(t)$ . Specifically, at the token  
 166 level, a token  $v \in \mathcal{V}$  can be deleted (denoted by transition to an empty state  $\emptyset$ ) with an infinitesimal rate  $\sigma(t)$ ,  
 167 or remain unchanged otherwise. Thus, the transition probability within infinitesimal time  $\Delta t$  is:

$$168 p_{t+\Delta t|t}(v'|v) = \begin{cases} \sigma(t)\Delta t, & v' = \emptyset, \\ 1 - \sigma(t)\Delta t, & v' = v. \end{cases} \quad (4)$$

169 Based on this independent token-level process, the sequence-level transition probability between timesteps  $s$   
 170 and  $t$  with  $0 < s < t < 1$  can be derived as:

$$171 p_{t|s}(\mathbf{x}_t|\mathbf{x}_s) = (1 - e^{-(\bar{\sigma}(t)-\bar{\sigma}(s))|\mathbf{x}_s|-|\mathbf{x}_t|})e^{-(\bar{\sigma}(t)-\bar{\sigma}(s))|\mathbf{x}_t|}N(\mathbf{x}_t, \mathbf{x}_s). \quad (5)$$

172 Here,  $|\mathbf{x}|$  denotes the length of the sequence  $\mathbf{x}$ ,  $\bar{\sigma}(t) = \int_0^t \sigma(\tau)d\tau$  is the integrated noise rate, and  $N(\mathbf{x}_t, \mathbf{x}_s)$   
 173 is the number of occurrences of  $\mathbf{x}_t$  as distinct subsequences in  $\mathbf{x}_s$ .<sup>1</sup> This number accounts for the multiplicity  
 174 of all the possible independent deletion paths from  $\mathbf{x}_s$  to  $\mathbf{x}_t$ ; see Appendix D.2 for the proof of Eq. 5.

175 The infinitesimal transition of the forward process is captured by the sequence-level transition rate matrix  $Q_t$ .  
 176 Due to token-wise independence, at most one deletion can occur within an infinitesimal time interval with  
 177 non-negligible ( $\Omega(\Delta t)$ ) probability. Thus, the rate  $Q_t(\mathbf{y}, \mathbf{x}_t)$  is non-zero only when  $\mathbf{y} = \mathbf{x}_t$  or  $\mathbf{y} \succ_1 \mathbf{x}_t$ , i.e.  
 178  $\mathbf{x}_t$  is the result of deleting exactly one token from  $\mathbf{y}$ , and the rate for  $\mathbf{y} \succ_1 \mathbf{x}_t$  is (details in Appendix D.3):

$$179 Q_t(\mathbf{y}, \mathbf{x}_t) = \lim_{\Delta t \rightarrow 0} \frac{p_{t+\Delta t|t}(\mathbf{x}_t|\mathbf{y})}{\Delta t} = \sigma(t)N(\mathbf{x}_t, \mathbf{y}). \quad (6)$$

180 <sup>1</sup>For example,  $N(\langle \text{BOS} \rangle \mathbf{b} \mathbf{a} \mathbf{g}, \langle \text{BOS} \rangle \mathbf{b} \mathbf{a} \mathbf{b} \mathbf{g} \mathbf{b} \mathbf{a} \mathbf{g}) = 5$ . The distinct subsequences are (highlighted in bold):  
 181  $\langle \text{BOS} \rangle \mathbf{b} \mathbf{a} \mathbf{b} \mathbf{g} \mathbf{b} \mathbf{a} \mathbf{g}$ ,  $\langle \text{BOS} \rangle \mathbf{b} \mathbf{a} \mathbf{b} \mathbf{g} \mathbf{b} \mathbf{a} \mathbf{g}$ ,  $\langle \text{BOS} \rangle \mathbf{b} \mathbf{a} \mathbf{b} \mathbf{g} \mathbf{b} \mathbf{a} \mathbf{g}$ ,  $\langle \text{BOS} \rangle \mathbf{b} \mathbf{a} \mathbf{b} \mathbf{g} \mathbf{b} \mathbf{a} \mathbf{g}$ ,  $\langle \text{BOS} \rangle \mathbf{b} \mathbf{a} \mathbf{b} \mathbf{g} \mathbf{b} \mathbf{a} \mathbf{g}$ .

Note that, in our implementation, we prepend an undeletable special token  $\langle \text{BOS} \rangle$  at the beginning of each sequence, so the above derivations should exclude  $\langle \text{BOS} \rangle$ . Therefore, the fully noised sequence is a single  $\langle \text{BOS} \rangle$ , which also serves as the initial input token at the first generation step to represent an empty sequence.

### 3.2 BACKWARD PROCESS: INSERTION

As in Sec. 2.1, we aim to learn the time reversal of the forward process, a CTMC with rate matrix  $\tilde{Q}_t(\mathbf{x}_t, \mathbf{y}) = Q_t(\mathbf{y}, \mathbf{x}_t) s(\mathbf{x}_t, t)_{\mathbf{y}}$ , where  $s$  is the **concrete score**. Since  $Q_t$  involves single-token deletions, the backward process  $\tilde{Q}_t$  only considers single-token insertions (i.e.,  $\mathbf{y} \succ_1 \mathbf{x}_t$ ).

Directly applying the concrete score learning approach used in MDLMs is impractical here. The reason is that given  $\mathbf{x}_t$ , the number of possible resulting states  $\mathbf{y}$  is variable, because different insertions can lead to the same result.<sup>2</sup> Consequently, targeting the concrete score requires calculating the number of possible  $\mathbf{y}$  values and enabling the model to produce a variable-shaped output representing the concrete score for each  $\mathbf{y}$ . These requirements collectively introduce significant complexity and implementation challenges.

As an alternative, we target the **insertion score**  $\bar{s}$ . We learn a score for every insertion, regardless of whether the resulting  $\mathbf{y}$  are identical. We define an insertion action  $(i, v)$  as inserting token  $v$  **after**<sup>3</sup> the  $i$ -th position of  $\mathbf{x}_t$ , resulting in  $\text{Ins}(\mathbf{x}_t, i, v) = (\mathbf{x}_{\leq i}, v, \mathbf{x}_{> i})$ . The definition of insertion score is:

$$\bar{s}(\mathbf{x}_t, t)[i, v] \stackrel{\text{def}}{=} \frac{\mathbb{E}_{\mathbf{x}_0}[(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0|} N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)]}{\mathbb{E}_{\mathbf{x}_0}[(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0|} N(\mathbf{x}_t, \mathbf{x}_0)]}, \quad \forall (i, v) \in [0, |\mathbf{x}_t|]_{\mathbb{Z}} \times \mathcal{V}, \quad (7)$$

whose shape is  $|\mathbf{x}_t| \times |\mathcal{V}|$ , tractable for transformer-based models.

Since the CTMC dynamics are based on the concrete score, in order to sample based on the insertion score, we first show that the concrete score is an average of insertion scores, and the reverse transition rate is a summation of the rate of insertion actions (details in Appendix D.4):

$$s(\mathbf{x}_t, t)_{\mathbf{y}} = \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \frac{1}{N(\mathbf{x}_t, \mathbf{y})} \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} \bar{s}(\mathbf{x}_t, t)[i, v(\mathbf{x}_t, \mathbf{y})], \quad (8)$$

$$\tilde{Q}_t(\mathbf{x}_t, \mathbf{y}) = \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} \underbrace{\left( \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \bar{s}(\mathbf{x}_t, t)[i, v(\mathbf{x}_t, \mathbf{y})] \right)}_{\text{Rate of action } (i, v(\mathbf{x}_t, \mathbf{y}))}, \quad (9)$$

where  $I(\mathbf{x}_t, \mathbf{y})$  is the set of viable insertion positions and  $v(\mathbf{x}_t, \mathbf{y})$  is the inserted token from  $\mathbf{x}_t$  to  $\mathbf{y}$ .<sup>2</sup>

Based on Eq. 9, we can transform the concrete score-based sampling in Eq. 2 into an equivalent insertion score-based sampling without computing  $N(\mathbf{x}_t, \mathbf{y})$ ,  $I(\mathbf{x}_t, \mathbf{y})$ , or  $v(\mathbf{x}_t, \mathbf{y})$  (details in Appendix D.5):

$$p_{t-\Delta t}^{\theta}((i, v)|\mathbf{x}_t) = \begin{cases} \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \bar{s}_{\theta}(\mathbf{x}_t, t)[i, v] \Delta t, & v \neq \emptyset, \\ 1 - \sum_{w \neq \emptyset} p_{t-\Delta t}^{\theta}((i, w)|\mathbf{x}_t), & v = \emptyset, \end{cases} \quad (10)$$

where  $\emptyset$  indicates no insertion, and Tau-leaping (Gillespie, 2001), a very popular approximate simulation method, could be adopted to sample all insertions simultaneously for parallel decoding.

### 3.3 TRAINING OBJECTIVE: DENOISING INSERTION SCORE ENTROPY

We aim to train the insertion score  $\bar{s}_{\theta}$  using the DSE objective (Eq. 3) with the derived transition rate  $Q(\mathbf{y}, \mathbf{x}_t)$  (Eq. 6), conditional distribution  $p_{t|s}(\mathbf{x}_t|\mathbf{x}_s)$  (Eq. 5), and parameterized concrete score  $s_{\theta}(\mathbf{x}_t, t)_{\mathbf{y}}$  (Eq. 8) for

<sup>2</sup>For example, inserting ‘a’ after the 1st or 2nd index of ‘ $\langle \text{BOS} \rangle$  b a g’ both yield ‘ $\langle \text{BOS} \rangle$  b a a g’.

<sup>3</sup>We use a prepended, non-deletable  $\langle \text{BOS} \rangle$  token (index 0) for insertions at the start.

235 DID. Here, directly substituting  $s_\theta$  (Eq. 8) into DSE is challenging. Since  $s_\theta$  is an average of insertion scores  
 236  $\bar{s}_\theta$ , it results in an intractable log-sum structure within the DSE objective (Appendix D.6), we apply Jensen’s  
 237 inequality to derive a tractable variational upper bound, the Denoising Insertion Score Entropy (DISE).

238 **Proposition 1** (Denoising Insertion Score Entropy (DISE)). The DSE objective for the deletion-insertion  
 239 process is upper bounded by the DISE objective,  $\mathcal{L}_\theta^{\text{DSE}}(\mathbf{x}_0) \leq \mathcal{L}_\theta^{\text{DISE}}(\mathbf{x}_0)$ , which is defined as:

$$241 \mathcal{L}_\theta^{\text{DISE}}(\mathbf{x}_0) = \mathbb{E}_{t, \mathbf{x}_t} \left\{ \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{i,v} \left[ \bar{s}_\theta(\mathbf{x}_t, t)[i, v] - \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \log \bar{s}_\theta(\mathbf{x}_t, t)[i, v] + C \right] \right\}, \quad (11)$$

242 where  $C = K \left( \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \right)$  is a  $\theta$ -free constant,  $t \sim \text{Unif}([0, 1])$ , and  $\mathbf{x}_t \sim p_{t|0}(\mathbf{x}_t | \mathbf{x}_0)$ .

243 The proof (Appendix D.6) utilizes Jensen’s inequality and a summation identity (Lemma 1) to transform the  
 244 objective from state-based ( $\mathbf{y}$ ) to action-based ( $(i, v)$ ).

### 245 3.4 EFFICIENT PARALLEL DYNAMIC PROGRAMMING FOR SUBSEQUENCE COUNTING PROBLEMS

246 A fundamental challenge for the DISE objective (Eq. 11) is to efficiently solve the ratios of  $N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)$   
 247 and  $N(\mathbf{x}_t, \mathbf{x}_0)$  for all possible insertions of  $(i, v)$  on  $\mathbf{x}_t$ . Suppose the lengths of  $\mathbf{x}_0$  and  $\mathbf{x}_t$  are  $m$  and  $n$ ,  
 248 solving a single subsequence counting problem of  $N(\mathbf{x}_t, \mathbf{x}_0)$  has a well-known time complexity of  $O(mn)$   
 249 through dynamic programming, performing it naively for  $n \times V$  times would be prohibitive. Here, we  
 250 show that  $N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)$  for all  $(i, v)$  pairs could be efficiently solved based on the intermediate results  
 251 of solving  $N(\mathbf{x}_t, \mathbf{x}_0)$  just twice (via a prefix DP in Eq. 13 and a suffix DP in Eq. 14), reducing the time  
 252 complexity to compute all ratios from  $O(mn^2V)$  to  $O(mn)$ , and making the training of DID possible.

253 **Counting**  $N(\mathbf{x}_t, \mathbf{x}_0)$ . Here we briefly introduce the classic prefix DP and suffix DP solutions to this problem.  
 254 Using Python’s slicing syntax, the base cases of empty sequences are:

$$255 N(\mathbf{x}_t[:0], \mathbf{x}_0[:j]) = N(\mathbf{x}_t[n:], \mathbf{x}_0[j:]) = 1, \quad \forall j \in \{0, \dots, m\}. \quad (12)$$

256 The **prefix DP** iteratively computes  $N(\mathbf{x}_t[:i], \mathbf{x}_0[:j])$  from the solved subproblems  $N(\mathbf{x}_t[:i], \mathbf{x}_0[:j-1])$   
 257 and  $N(\mathbf{x}_t[:i-1], \mathbf{x}_0[:j-1])$  to get the final result  $N(\mathbf{x}_t[:n], \mathbf{x}_0[:m])$ , and the state transition is:

$$258 N(\mathbf{x}_t[:i], \mathbf{x}_0[:j]) = N(\mathbf{x}_t[:i], \mathbf{x}_0[:j-1]) + \delta(\mathbf{x}_t[i-1], \mathbf{x}_0[j-1]) \cdot N(\mathbf{x}_t[:i-1], \mathbf{x}_0[:j-1]). \quad (13)$$

259 The **suffix DP** iteratively computes  $N(\mathbf{x}_t[i:], \mathbf{x}_0[j:])$  from the solved subproblems  $N(\mathbf{x}_t[i:], \mathbf{x}_0[j+1:])$   
 260 and  $N(\mathbf{x}_t[i+1:], \mathbf{x}_0[j+1:])$  to get the final result  $N(\mathbf{x}_t[0:], \mathbf{x}_0[0:])$ , and the state transition is:

$$261 N(\mathbf{x}_t[i:], \mathbf{x}_0[j:]) = N(\mathbf{x}_t[i:], \mathbf{x}_0[j+1:]) + \delta(\mathbf{x}_t[i], \mathbf{x}_0[j]) \cdot N(\mathbf{x}_t[i+1:], \mathbf{x}_0[j+1:]). \quad (14)$$

262 The time complexities are  $O(mn)$  for both the prefix and suffix DPs. Notably, they could be batched and  
 263 parallelized along the  $i$ -dimension, thus supporting vectorization and parallel execution, and it only needs to  
 264 sequentially loop over the  $j$ -dimension for  $m$  times.

265 **Counting**  $N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)$ . It could be efficiently solved based on the results of prefix and suffix DP:

$$266 N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0) = \underbrace{\sum_{j=1}^m \left[ \delta(\mathbf{x}_0[j], v) \cdot N(\mathbf{x}_t[:i], \mathbf{x}_0[:j-1]) \cdot N(\mathbf{x}_t[i:], \mathbf{x}_0[j:]) \right]}_{\text{index addition}}, \quad (15)$$

267 due to the form of Eq. 15, results for all  $(i, v)$  pairs can be solved in parallel with an elementwise multiplication  
 268 of the prefix and suffix DP result matrices, followed by an index addition that could be efficiently implemented  
 269 with a sparse tensor coalescence. A PyTorch implementation of the DP algorithms is in Appendix G.4.

### 3.5 SIMPLIFIED MODEL FOR FIXED-LENGTH SETTING

To facilitate a fair comparison with MDLMs on the widely-adopted fixed-length language modeling benchmarks (Tab. 1), and clearly isolate the superior FLOPs efficiency of DID, we develop a set of optimizations to enhance DID in the fixed-length setting. We show that when  $|\mathbf{x}_0|$  is a constant, 1) the insertion score becomes time-independent as the time-dependent terms of  $(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0|}$  in Eq. 7 could be canceled out, which leads to 2) a sequence-level normalization property (details in Appendix D.7):

$$\sum_{i,v} \bar{s}(\mathbf{x}_t, t)[i, v] = |\mathbf{x}_0| - |\mathbf{x}_t|. \quad (16)$$

This benefits the parameterization and training of DID from two aspects. First, the time-independence means that the network does not require time  $t$  as input in the fixed-length setting, thus the parameterization reduces to  $\bar{s}_\theta(\mathbf{x}_t)$ , saves the parameters for time embedding, and enables a cache mechanism similar to (Ou et al., 2025) if the sequence is not changed between steps. Second, the output of the insertion score network could be explicitly normalized with a summation of  $|\mathbf{x}_0| - |\mathbf{x}_t|$  as in Eq. 16. Therefore, the summation term of outputs from the insertion score network in the DISE objective (Eq. 11) turns into a constant, giving rise to a simplified Denoising Insertion Cross Entropy (DICE) objective (details in Appendix D.8):

$$\mathcal{L}_\theta^{\text{DICE}}(\mathbf{x}_0) = \mathbb{E}_{t, \mathbf{x}_t} \left\{ \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{i,v} \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \left[ -\log \bar{s}_\theta(\mathbf{x}_t)[i, v] + C \right] \right\}, \quad (17)$$

where  $C = \log \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)}$  is a  $\theta$ -free constant. DICE can be interpreted as a weighted cross-entropy loss between the predicted insertion scores and the ground truth of subsequence count ratios, hence the name.

## 4 EXPERIMENTS

We evaluate language modeling performance, sampling quality, and training/inference efficiency of DID in the fixed-length setting in Sec. 4.1 and the variable-length setting in Sec. 4.2. For more details, results, generation examples, and intermediate generation process, please refer to Appendix E, F, H.

### 4.1 DID FOR FIXED-LENGTH LANGUAGE MODELING

**Settings.** Following RADD (Ou et al., 2025), which serves as our baseline method of MDLM, we train DID of both small and medium sizes on the OpenWebText (OWT) dataset (Gokaslan & Cohen, 2019) with the DICE objective (Eq. 17). We adopt the GPT2 tokenizer (Radford et al., 2019), concatenate all sequences, and split them into fixed-length chunks of 1024 tokens, and the training batch size is 512. RADD-small and -medium are reproduced with their open-sourced model checkpoints (trained for 400K steps on OWT).

**Zero-shot language modeling perplexity.** We evaluate the zero-shot modeling perplexity on seven datasets in Tab. 1. Since DID eliminates the computational FLOPs for  $\langle \text{MASK} \rangle$  and hence reduces FLOPs by approximately half compared to RADD for the same training steps, we compare DID under two training configurations: **DID-S** (Steps-aligned), trained for the same steps as RADD (400K) and **DID-F** (FLOPs-aligned), trained for double the steps (800K) to match the total computational budget. We observe that when aligned by training steps (DID-S), our model achieves performance comparable to RADD, despite utilizing only about half the computational FLOPs. When aligned by the total computational budget (DID-F), DID consistently outperforms RADD across the majority of datasets for both small and medium sizes. This demonstrates that the computational savings achieved by eliminating  $\langle \text{MASK} \rangle$  tokens are effectively turned into improved modeling performance. We also provide an ablation study in Appendix F.1 for DID trained with DISE objective (Eq. 11), i.e., without the additional optimizations in DICE for fixed-length data introduced in Sec. 3.5. This results in a reduced performance than DICE-trained DID in Tab. 1, yet comparable to RADD.

Table 1: Zero-shot language modeling perplexity. Results for diffusion models are perplexity upper bounds.

Size	Method	WikiText	Lambada	Pubmed	AG News	LM1B	Arxiv	PTB
Small	RADD	<u>38.27</u>	51.82	56.99	<u>73.18</u>	<u>72.99</u>	85.95	<b>108.79</b>
	DID-S	<u>38.72</u>	<u>49.10</u>	<u>55.02</u>	<u>76.02</u>	<u>74.04</u>	<u>82.41</u>	115.37
	DID-F	<b>36.91</b>	<b>48.00</b>	<b>52.89</b>	<b>71.48</b>	<b>72.04</b>	<b>78.38</b>	<u>111.60</u>
Medium	RADD	<u>28.44</u>	44.10	41.06	<u>48.96</u>	60.32	66.28	<b>81.05</b>
	DID-S	29.19	<u>41.94</u>	<u>40.84</u>	52.53	<u>59.88</u>	<u>63.95</u>	91.87
	DID-F	<b>28.35</b>	<b>41.00</b>	<b>38.71</b>	<b>48.84</b>	<b>58.05</b>	<b>61.77</b>	<u>87.09</u>

Table 2: Generative perplexity (PPL, evaluated by GPT2 Large), unigram entropy, inference time (in seconds), speedup, and average generation length for fixed-length models under different total denoising steps.

Method	Steps	16	32	64	128	256	512	1024
RADD	PPL	284.78	155.01	111.56	95.10	87.56	<b>84.00</b>	<b>84.05</b>
	Entropy	8.35	8.26	8.20	8.15	8.11	8.10	8.09
	Time (s)	0.220	0.317	0.499	0.879	1.644	2.882	4.512
DID	PPL	<b>158.93</b>	<b>110.06</b>	<b>97.32</b>	<b>91.25</b>	<b>86.98</b>	86.04	85.35
	Entropy	8.15	8.13	8.13	8.12	8.09	8.08	8.09
	Time (s)	<b>0.169</b>	<b>0.246</b>	<b>0.353</b>	<b>0.573</b>	<b>1.047</b>	<b>1.826</b>	<b>3.006</b>
	Speedup	1.30×	1.29×	1.41×	1.53×	1.57×	1.58×	1.50×
	Length	1023.29	1024.01	1024.18	1024.07	1023.91	1024.03	1024.10

**Generative performance.** We report generative perplexity (PPL), unigram entropy (diversity), and inference speed of direct sampling across different numbers of total denoising steps in Tab. 2. Compared with RADD, DID achieves significantly better generation quality (lower PPL) with fewer denoising steps. When more total denoising steps are used, RADD slightly outperforms DID, which is reasonable since RADD is naturally designed for the fixed-length setting. Nonetheless, DID could consistently outperform RADD in the variable-length setting, which is deferred to Tab. 4. Furthermore, DID consistently provides  $\sim 1.5\times$  inference speedup. This improvement stems from the fact that the average sequence length during the iterative insertion process is shorter than the fixed-length process of RADD. We also provide nucleus sampling results in Appendix F.3, where DID achieves lower PPL and entropy compared with RADD, demonstrating a stronger annealing effect.

**Training speed.**<sup>4</sup> We report the training speeds for different model sizes in Tab. 3 to verify the efficiency gains of DID. DID demonstrates substantial speedups, increasing from  $1.63\times$  to  $1.82\times$  as we scale up the model. Empirically, removing the computations for `<MASK>` tokens significantly boosts training efficiency.

While the reduction in FLOPs suggests a theoretical  $2\times$  speedup, the actual gains in Tab. 2 are more modest. This discrepancy, as discussed in (Zheng et al., 2025), arises because inference is not a purely FLOPs-bound task. For training, the observed gains in Tab. 3 are also slightly lower due to additional overheads such as the DP algorithm for loss implementation (Sec. 3.4), which is a constant overhead independent of the model size, thus the speedup could be scaled up with model size. Besides, a less developed system-level support for variable-length data also constrains the speedups.

Table 3: Average training time (in seconds) per 50 steps (i.e. batches) on OpenWebText.

	Small	Medium	Large
RADD	26.46	53.17	92.90
DID	<b>16.26</b>	<b>31.16</b>	<b>51.14</b>
Speedup	1.63×	1.71×	1.82×

## 4.2 DID FOR VARIABLE-LENGTH LANGUAGE MODELING

<sup>4</sup>Large models are not fully trained; only their training speeds are measured.



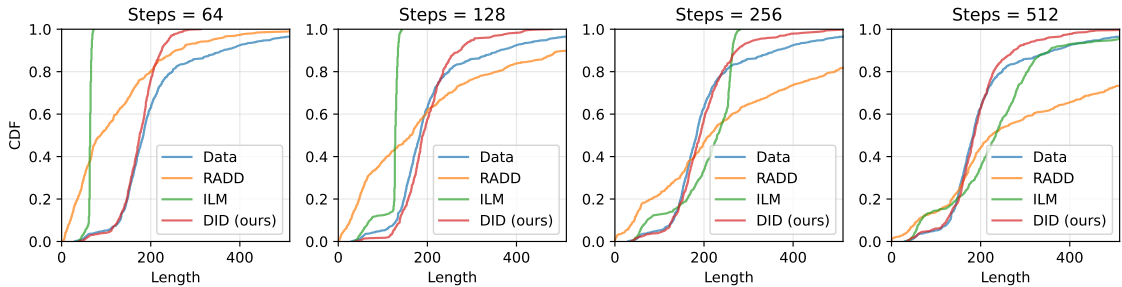


Figure 2: Cumulative distribution functions (CDFs) of generation length under different total denoising steps.

**Settings.** Following ILM (Patel et al., 2025), an insertion-based approach, we train DID-small on the Stories dataset (Eldan & Li, 2023; Mostafazadeh et al., 2016) for 60K steps with batch size 512, utilizing its variable-length sequences (average length 213.43 under the Bert-base-uncased tokenizer (Devlin et al., 2018)) truncated to a maximum length of 1024 and without padding. We also train RADD-small on Stories for the same steps ( $> 2 \times$  FLOPs of DID), details in Appendix E. Since RADD requires fixed-length inputs, we pad all sequences to a length of 1024, which is in line with the setting for MDLM in ILM experiments (Patel et al., 2025). As a result, RADD generates fixed-length outputs containing  $\langle \text{PAD} \rangle$  tokens, which we subsequently remove to obtain the final variable-length outputs. ILM is reproduced with its open-sourced checkpoints.

**Generative performance.** We report the generation quality and speed of direct sampling for variable-length models in Tab. 4. Compared with both baselines, DID maintains a significantly lower generative PPL (evaluated by GPT2 Large) with relatively high diversity and is more stable across different numbers of steps. This is important because strong sensitivity and non-convergence to the manually predefined number of generation steps are undesirable.

Regarding inference speed, DID achieves a speedup of up to  $3.79 \times$  compared to RADD, due to the savings of  $\langle \text{MASK} \rangle$  and  $\langle \text{PAD} \rangle$  tokens for both model calling and distribution sampling. On the other hand, ILM achieves the fastest speed despite its lowest quality. We credit this to its over-simplified generation mechanism, which samples a token at exactly one *designated* position per step. In contrast, DID and RADD sample from all possible token positions to enable parallel decoding. ILM benefits considerably from this simplification, since categorical sampling is a major bottleneck in small models (Zheng et al., 2025). Moreover, ILM is limited to generating text shorter than the total steps (see Tab. 4), which also contributes to its faster sampling.

We also provide nucleus sampling results for variable-length models in Appendix F.3, and ablation studies of different padding lengths for RADD in Appendix F.4, addressing potential concerns that the 1024 padding length for RADD might be too long or unfair. When trained at a length of 512, RADD exhibits observable degradation and remains  $\sim 1.59 \times$  slower than DID, further confirming the original setting of ILM.

Table 4: Generative PPL, unigram entropy, inference time (in seconds), and average generation length for variable-length models under different denoising steps. \*: as outliers significantly affect PPL, only samples with PPL  $< 300$  are counted,  $\dagger$ : speedup over RADD.

Method	Steps	64	128	256	512
ILM	PPL*	161.80	137.64	42.29	31.14
	Entropy	5.20	5.65	5.97	6.01
	Time (s)	<b>0.016</b>	<b>0.034</b>	<b>0.087</b>	<b>0.271</b>
	Length	63.34	120.77	206.44	234.44
RADD	PPL*	<u>81.92</u>	<u>50.89</u>	<u>34.47</u>	<u>26.78</u>
	Entropy	5.22	5.58	5.79	5.85
	Time (s)	0.246	0.441	0.827	1.461
	Length	110.66	200.73	349.54	353.47
DID	PPL	<b>22.78</b>	<b>21.07</b>	<b>21.90</b>	<b>23.88</b>
	Entropy	5.90	5.94	5.94	5.94
	Time (s)	<u>0.090</u>	<u>0.132</u>	<u>0.218</u>	<u>0.388</u>
	Speedup $\dagger$	$2.73 \times$	$3.34 \times$	$3.79 \times$	$3.76 \times$
	Length	182.31	193.77	202.97	204.96

**Length modeling.** Besides, DID demonstrates superior length modeling capabilities, exhibiting consistency between the generation length distribution and the training data length distribution. This is demonstrated in Fig. 2, where the CDF of the length distribution of DID is closely aligned with the dataset compared to the baselines. DID’s average generation length reported in Tab. 4 is also stable and approximating the ground truth distribution.

**Training speed.**<sup>5</sup> We also compare the training speed on the Stories dataset (Tab. 5). The efficiency gains of DID are even more pronounced in the variable-length setting, reaching up to  $3.10\times$  speedup for large models. This is because RADD suffers significant overhead from processing  $\langle\text{PAD}\rangle$  tokens (since the average length 213.43 is much shorter than the padded length 1024), while DID benefits from the shorter length.

Regarding the language modeling ability for variable-length models, please see Appendix F.2.

## 5 CONCLUSION

In this paper, we introduce DID to improve the computational efficiency and generation flexibility of diffusion language models by eliminating the use of  $\langle\text{MASK}\rangle$  and  $\langle\text{PAD}\rangle$  tokens in our paradigm. Theoretically, we formulate the diffusion processes for deletion and insertion, define an insertion score, derive and implement the corresponding training objectives and sampling algorithm for DID. We evaluated DID on modeling performance, generation quality, and training/inference speed, demonstrating the superiority of DID over the baselines of MDLM and existing insertion-based LMs in both fixed-length and variable-length settings. A discussion of the limitations and future works of this paper is provided in Appendix A.

Table 5: Average training time (in seconds) per 50 steps on Stories.

	Small	Medium	Large
RADD	19.93	37.87	67.75
DID	<b>10.59</b>	<b>14.20</b>	<b>21.83</b>
Speedup	$1.88\times$	$2.67\times$	$3.10\times$

<sup>5</sup>Medium and large models are not fully trained; only their training speeds are measured.

## REFERENCES

- Marianne Arriola, Aaron Gokaslan, Justin T Chiu, Zhihan Yang, Zhixuan Qi, Jiaqi Han, Subham Sekhar Sahoo, and Volodymyr Kuleshov. Block diffusion: Interpolating between autoregressive and diffusion language models. *arXiv preprint arXiv:2503.09573*, 2025.
- Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in neural information processing systems*, 34:17981–17993, 2021.
- Andrew Campbell, Joe Benton, Valentin De Bortoli, Thomas Rainforth, George Deligiannidis, and Arnaud Doucet. A continuous time framework for discrete denoising models. *Advances in Neural Information Processing Systems*, 35:28266–28279, 2022.
- Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and Tony Robinson. One billion word benchmark for measuring progress in statistical language modeling, 2014. URL <https://arxiv.org/abs/1312.3005>.
- Arman Cohan, Franck Dernoncourt, Doo Soon Kim, Trung Bui, Seokhwan Kim, Walter Chang, and Nazli Goharian. A discourse-aware attention model for abstractive summarization of long documents. *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, 2018. doi: 10.18653/v1/n18-2097. URL <http://dx.doi.org/10.18653/v1/n18-2097>.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. In *Advances in Neural Information Processing Systems*, 2022.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018. URL <http://arxiv.org/abs/1810.04805>.
- Ronen Eldan and Yuanzhi Li. Tinstories: How small can language models be and still speak coherent english?, 2023. URL <https://arxiv.org/abs/2305.07759>.
- Daniel T Gillespie. Approximate accelerated stochastic simulation of chemically reacting systems. *The Journal of chemical physics*, 115(4):1716–1733, 2001.
- Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. <http://Skylion007.github.io/OpenWebTextCorpus>, 2019.
- Shansan Gong, Shivam Agarwal, Yizhe Zhang, Jiacheng Ye, Lin Zheng, Mukai Li, Chenxin An, Peilin Zhao, Wei Bi, Jiawei Han, et al. Scaling diffusion language models via adaptation from autoregressive models. *arXiv preprint arXiv:2410.17891*, 2024.
- Jiatao Gu, Changan Wang, and Junbo Zhao. Levenshtein transformer. In *Advances in Neural Information Processing Systems*, volume 32, pp. 11179–11189, 2019.
- Marton Havasi, Brian Karrer, Itai Gat, and Ricky TQ Chen. Edit flows: Flow matching with edit operations. *arXiv preprint arXiv:2506.09018*, 2025.
- Jaeyeon Kim, Lee Cheuk-Kit, Carles Domingo-Enrich, Yilun Du, Sham Kakade, Timothy Ngotiaoco, Sitan Chen, and Michael Albergo. Any-order flexible length masked diffusion, 2025. URL <https://arxiv.org/abs/2509.01025>.

- 517 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on*  
518 *Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.  
519
- 520 Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete Diffusion Modeling by Estimating the Ratios of the  
521 Data Distribution. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235,  
522 pp. 32819–32848. PMLR, 2024.
- 523 Sidi Lu, Tao Meng, and Nanyun Peng. Insnet: An efficient, flexible, and performant insertion-based text  
524 generation model. In *Advances in Neural Information Processing Systems*, volume 35, 2022.  
525
- 526 Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a large annotated corpus of  
527 English: The Penn Treebank. *Computational Linguistics*, 19(2):313–330, 1993. URL <https://www.aclweb.org/anthology/J93-2004>.  
528
- 529 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. In  
530 *International Conference on Learning Representations*, 2017.  
531
- 532 Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende,  
533 Pushmeet Kohli, and James Allen. A corpus and cloze evaluation for deeper understanding of commonsense  
534 stories. In Kevin Knight, Ani Nenkova, and Owen Rambow (eds.), *Proceedings of the 2016 Conference*  
535 *of the North American Chapter of the Association for Computational Linguistics: Human Language*  
536 *Technologies*, pp. 839–849, San Diego, California, June 2016. Association for Computational Linguistics.  
537 doi: 10.18653/v1/N16-1098. URL <https://aclanthology.org/N16-1098/>.
- 538 Shen Nie, Fengqi Zhu, Chao Du, Tianyu Pang, Qian Liu, Guangtao Zeng, Min Lin, and Chongxuan Li.  
539 Scaling up masked diffusion models on text. *arXiv preprint arXiv:2410.18514*, 2024.  
540
- 541 Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-Rong Wen,  
542 and Chongxuan Li. Large language diffusion models. *arXiv preprint arXiv:2502.09992*, 2025.  
543
- 544 Manfred Opper and Guido Sanguinetti. Variational inference for markov jump processes. In J. Platt, D. Koller,  
545 Y. Singer, and S. Roweis (eds.), *Advances in Neural Information Processing Systems*, volume 20. Curran  
546 Associates, Inc., 2007. URL [https://proceedings.neurips.cc/paper\\_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2007/file/735b90b4568125ed6c3f678819b6e058-Paper.pdf)  
547 [2007/file/735b90b4568125ed6c3f678819b6e058-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2007/file/735b90b4568125ed6c3f678819b6e058-Paper.pdf).
- 548 Jingyang Ou, Shen Nie, Kaiwen Xue, Fengqi Zhu, Jiacheng Sun, Zhenguo Li, and Chongxuan Li. Your  
549 absorbing discrete diffusion secretly models the conditional distributions of clean data. In *The Thirteenth*  
550 *International Conference on Learning Representations*, 2025. URL [https://openreview.net/](https://openreview.net/forum?id=7yqjVgWWxx)  
551 [forum?id=7yqjVgWWxx](https://openreview.net/forum?id=7yqjVgWWxx).
- 552 Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc Quan Pham, Raffaella Bernardi, Sandro  
553 Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernandez. The LAMBADA dataset: Word prediction  
554 requiring a broad discourse context. In *Proceedings of the 54th Annual Meeting of the Association*  
555 *for Computational Linguistics (Volume 1: Long Papers)*, pp. 1525–1534, Berlin, Germany, August  
556 2016. Association for Computational Linguistics. URL [http://www.aclweb.org/anthology/](http://www.aclweb.org/anthology/P16-1144)  
557 [P16-1144](http://www.aclweb.org/anthology/P16-1144).
- 558 Dhruv Patel, Aishwarya Sahoo, Avinash Amballa, Tahira Naseem, Tim G. J. Rudner, and Andrew  
559 McCallum. Insertion language models: Sequence generation with arbitrary-position insertions, 2025. URL  
560 <https://arxiv.org/abs/2505.05755>.  
561
- 562 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the*  
563 *IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.

- 564 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models  
565 are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 566
- 567 Subham Sahoo, Marianne Arriola, Yair Schiff, Aaron Gokaslan, Edgar Marroquin, Justin Chiu, Alexander  
568 Rush, and Volodymyr Kuleshov. Simple and effective masked diffusion language models. *Advances in*  
569 *Neural Information Processing Systems*, 37:130136–130184, 2024.
- 570 Subham Sekhar Sahoo, Zhihan Yang, Yash Akhauri, Johnna Liu, Deepansha Singh, Zhoujun Cheng,  
571 Zhengzhong Liu, Eric Xing, John Thickstun, and Arash Vahdat. Esoteric language models, 2025. URL  
572 <https://arxiv.org/abs/2506.01928>.
- 573
- 574 Jiaxin Shi, Kehang Han, Zhe Wang, Arnaud Doucet, and Michalis Titsias. Simplified and generalized masked  
575 diffusion for discrete data. *Advances in neural information processing systems*, 37:103131–103167, 2024.
- 576 Mitchell Stern, William Chan, Jamie Kiros, and Jakob Uszkoreit. Insertion transformer: Flexible sequence  
577 generation via insertion operations. In *Proceedings of the 36th International Conference on Machine*  
578 *Learning*, volume 97, pp. 5976–5985. PMLR, 2019.
- 579
- 580 Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced  
581 transformer with rotary position embedding, 2023. URL <https://arxiv.org/abs/2104.09864>.
- 582 Guanghan Wang, Yair Schiff, Subham Sahoo, and Volodymyr Kuleshov. Remasking discrete diffusion models  
583 with inference-time scaling. *arXiv preprint arXiv:2503.00307*, 2025.
- 584
- 585 Chengyue Wu, Hao Zhang, Shuchen Xue, Zhijian Liu, Shizhe Diao, Ligeng Zhu, Ping Luo, Song Han,  
586 and Enze Xie. Fast-dllm: Training-free acceleration of diffusion llm by enabling kv cache and parallel  
587 decoding, 2025a. URL <https://arxiv.org/abs/2505.22618>.
- 588 Zirui Wu, Lin Zheng, Zhihui Xie, Jiacheng Ye, Jiahui Gao, Yansong Feng, Zhenguo Li, Victoria W., Guorui  
589 Zhou, and Lingpeng Kong. Dreamon: Diffusion language models for code infilling beyond fixed-size  
590 canvas, 2025b. URL <https://hkunlp.github.io/blog/2025/dreamon>.
- 591
- 592 Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification.  
593 *Advances in neural information processing systems*, 28, 2015.
- 594 Kaiwen Zheng, Yongxin Chen, Hanzi Mao, Ming-Yu Liu, Jun Zhu, and Qinsheng Zhang. Masked diffusion  
595 models are secretly time-agnostic masked models and exploit inaccurate categorical sampling. In *The*  
596 *Thirteenth International Conference on Learning Representations*, 2025. URL [https://openreview.](https://openreview.net/forum?id=CTC7CmirNr)  
597 [net/forum?id=CTC7CmirNr](https://openreview.net/forum?id=CTC7CmirNr).
- 598
- 599
- 600
- 601
- 602
- 603
- 604
- 605
- 606
- 607
- 608
- 609
- 610

611	CONTENTS	
612		
613	<b>1 Introduction</b>	<b>1</b>
614		
615		
616	<b>2 Preliminaries and Related Works</b>	<b>2</b>
617	2.1 Continuous-Time Discrete Diffusion . . . . .	2
618	2.2 Insertion-Based Language Models . . . . .	3
619		
620		
621	<b>3 DID: Deletion-Insertion Diffusion Language Models</b>	<b>4</b>
622	3.1 Forward Process: Deletion . . . . .	4
623	3.2 Backward Process: Insertion . . . . .	5
624	3.3 Training Objective: Denoising Insertion Score Entropy . . . . .	5
625	3.4 Efficient Parallel Dynamic Programming for Subsequence Counting Problems . . . . .	6
626	3.5 Simplified Model for Fixed-Length Setting . . . . .	7
627		
628		
629		
630		
631	<b>4 Experiments</b>	<b>7</b>
632	4.1 DID for Fixed-Length Language Modeling . . . . .	7
633	4.2 DID for Variable-Length Language Modeling . . . . .	8
634		
635		
636	<b>5 Conclusion</b>	<b>10</b>
637		
638	<b>A Limitations and Future Works</b>	<b>16</b>
639		
640	<b>B The Use of Large Language Models (LLMs)</b>	<b>16</b>
641		
642	<b>C Notation Summary</b>	<b>16</b>
643		
644	<b>D Detailed Proofs and Derivations</b>	<b>17</b>
645	D.1 Derivation of Denoising Score Entropy Loss (Eq.3) . . . . .	17
646	D.2 Derivation of Sequence-Level Transition Probability (Eq.5) . . . . .	18
647	D.3 Derivation of the Transition Rate (Eq.6) . . . . .	19
648	D.4 Derivation of the Relationship between Concrete Score and Insertion Score (Eq. 8) . . . . .	19
649	D.5 Derivation of the Sampling Probability (Eq.10) . . . . .	20
650	D.6 Derivation of the DISE Objective (Eq.11) . . . . .	21
651	D.7 Derivation of the Sequence-Level Normalization Property (Eq. 16) . . . . .	24
652	D.8 Derivation of the DICE Objective for Fixed-Length Data (Eq. 17) . . . . .	25
653		
654		
655		
656		
657		

658	<b>E Experimental Details</b>	<b>26</b>
659	E.1 Fixed-Length Training Details . . . . .	26
660	E.2 Variable-Length Training Details . . . . .	27
661	E.3 Sampling Details . . . . .	28
662	E.4 Evaluation Metrics . . . . .	28
663		
664		
665		
666	<b>F More Experimental Results</b>	<b>28</b>
667	F.1 Ablation Study of Sequence-Level Normalization for Fixed-Length Models . . . . .	28
668	F.2 Language Modeling Performance for Variable-Length Models . . . . .	28
669	F.3 Nucleus Sampling Results . . . . .	29
670	F.4 Comparisons with RADD of Different Padding Lengths . . . . .	30
671		
672		
673		
674	<b>G Algorithm Details</b>	<b>32</b>
675	G.1 Dynamic Programming for Subsequence Counting . . . . .	32
676	G.2 DID Training Algorithm . . . . .	33
677	G.3 DID Inference Algorithm . . . . .	33
678	G.4 PyTorch Implementation of the Dynamic Programming Algorithms . . . . .	34
679		
680		
681		
682	<b>H Generation Examples</b>	<b>37</b>
683	H.1 Samples Generated by the Fixed-Length Model Trained on OpenWebText . . . . .	37
684	H.1.1 Unconditional Generation . . . . .	37
685	H.1.2 Conditional Generation . . . . .	44
686	H.2 Samples Generated by the Variable-Length Model Trained on Stories . . . . .	47
687	H.3 Demonstration of the Intermediate Generation Process . . . . .	52
688		
689		
690		
691		
692		
693		
694		
695		
696		
697		
698		
699		
700		
701		
702		
703		
704		

## A LIMITATIONS AND FUTURE WORKS

This work presents the core framework of DID and, unlike the more established MDLMs, has not yet integrated many optimizations, such as advanced inference algorithms (Wu et al., 2025a; Wang et al., 2025), or hybrid models combining autoregressive approaches (Arriola et al., 2025; Sahoo et al., 2025). Since these optimizations are not inherently tied to a specific diffusion process, adapting them to DID represents a promising future direction. Second, although we have demonstrated the effectiveness, efficiency, and flexibility of DID, our models were trained at a relatively small scale due to resource constraints. As a result, their performance on larger and more complex tasks remains unexplored, and we leave scaling up DID to future work.

## B THE USE OF LARGE LANGUAGE MODELS (LLMs)

Following ICLR guidelines, we wish to clarify our use of Large Language Models (LLMs) during the preparation of this work.

The research ideas, methodology, experimental design, and analysis presented in this paper were developed entirely by the human authors. LLMs were not involved in the ideation process.

We utilized LLMs as tools for editing and polishing the text, helping to improve the clarity and phrasing in various sections of the main paper and the appendix. The authors have reviewed the manuscript thoroughly and take full responsibility for its content.

## C NOTATION SUMMARY

Table 6: Summary of Key Notations

Notation	Description
$\mathcal{V}$	Vocabulary.
$\mathcal{X} = \bigcup_{d=0}^{\infty} \mathcal{V}^d$	Sequence state space (variable length).
$\mathbf{x}_0, \mathbf{x}_t, \mathbf{y}$	Clean sequence; sequence at time $t$ ; another sequence state.
$ \mathbf{x} $	Length of $\mathbf{x}$ .
$\langle \text{BOS} \rangle, \langle \text{MASK} \rangle, \langle \text{PAD} \rangle$	Begin-of-sequence (non-deletable), absorbing mask, padding.
$\emptyset$	Null token representing deletion or no insertion (not in $\mathcal{V}$ ).
$Q_t, \tilde{Q}_t$	Forward and reverse sequence-level CTMC rate matrices.
$p_{t s}(\cdot \cdot)$	Forward transition probability from $s$ to $t$ .
$p_t(\cdot)$	Marginal distribution at time $t$ .
$\sigma(t), \bar{\sigma}(t)$	Noise rate and its integral $\int_0^t \sigma(\tau) d\tau$ .
$s(\mathbf{x}_t, t)_{\mathbf{y}}$	Concrete score $p_t(\mathbf{y})/p_t(\mathbf{x}_t)$ .
$N(\mathbf{x}, \mathbf{y})$	The number of occurrences of $\mathbf{x}$ as a distinct subsequence of $\mathbf{y}$ .
$\mathbf{y} \succ_1 \mathbf{x}$	$\mathbf{y}$ is obtained by inserting exactly one token into $\mathbf{x}$ .
$v(\mathbf{x}, \mathbf{y})$	The unique inserted token when $\mathbf{y} \succ_1 \mathbf{x}$ .
$\text{Ins}(\mathbf{x}, i, v)$	The result of inserting token $v \in \mathcal{V}$ after position $i$ of $\mathbf{x}$ .
$I(\mathbf{x}, \mathbf{y})$	Valid insertion indices s.t. $\text{Ins}(\mathbf{x}, i, v(\mathbf{x}, \mathbf{y})) = \mathbf{y}$ .
$\bar{s}(\mathbf{x}_t, t)[i, v]$	Insertion score for insertion operation $(i, v)$ at time $t$ .
$\bar{s}(\mathbf{x}_t)[i, v]$	Time-independent insertion score (fixed-length setting).
$K(a)$	Convex function $a(\log a - 1)$ .



## D DETAILED PROOFS AND DERIVATIONS

### D.1 DERIVATION OF DENOISING SCORE ENTROPY LOSS (EQ.3)

The training objective for discrete diffusion models is derived by minimizing a variational upper bound on the negative log-likelihood (NLL) of the data,  $-\log p_0^\theta(\mathbf{x}_0)$ .

Let  $\mathbb{P}_{\mathbf{x}_0}$  denote the path measure (the probability distribution over entire trajectories) of the true posterior reverse process conditioned on the data  $\mathbf{x}_0$ . Let  $\mathbb{P}^\theta$  denote the path measure of the learned reverse process parameterized by  $\theta$ . By the data processing inequality, we have:

$$-\log p_0^\theta(\mathbf{x}_0) = D_{\text{KL}}(\delta_{\mathbf{x}_0} \| p_0^\theta) \leq D_{\text{KL}}(\mathbb{P}_{\mathbf{x}_0} \| \mathbb{P}^\theta). \quad (18)$$

We define the training objective as this variational upper bound:  $\mathcal{L}(\theta) = D_{\text{KL}}(\mathbb{P}_{\mathbf{x}_0} \| \mathbb{P}^\theta)$ , assuming both processes share the same prior distribution at  $t = 1$ .

Both processes are Continuous-Time Markov Chains (CTMCs). Let  $\tilde{Q}_t^0(\mathbf{x}_t, \mathbf{y})$  denote the true conditional reverse transition rate (for  $\mathbb{P}_{\mathbf{x}_0}$ ) and  $\tilde{Q}_t^\theta(\mathbf{x}_t, \mathbf{y})$  denote the parameterized reverse transition rate (for  $\mathbb{P}^\theta$ ).

The KL divergence between path measures can be decomposed using the chain rule. If we consider a discrete-time approximation with infinitesimal step  $\Delta t$ , the total KL divergence is the sum of the expected KL divergences at each step. We analyze the KL divergence between the infinitesimal transition probabilities  $p^0(\mathbf{y}|\mathbf{x}_t) = \delta(\mathbf{x}_t, \mathbf{y}) + \tilde{Q}_t^0(\mathbf{x}_t, \mathbf{y})\Delta t + o(\Delta t)$  and  $p^\theta(\mathbf{y}|\mathbf{x}_t) = \delta(\mathbf{x}_t, \mathbf{y}) + \tilde{Q}_t^\theta(\mathbf{x}_t, \mathbf{y})\Delta t + o(\Delta t)$ . The instantaneous KL divergence is (Oppor & Sanguinetti, 2007):

$$D_{\text{KL}}(p^0(\cdot|\mathbf{x}_t) \| p^\theta(\cdot|\mathbf{x}_t)) = \Delta t \sum_{\mathbf{y} \neq \mathbf{x}_t} \left( \tilde{Q}_t^0(\mathbf{x}_t, \mathbf{y}) \log \frac{\tilde{Q}_t^0(\mathbf{x}_t, \mathbf{y})}{\tilde{Q}_t^\theta(\mathbf{x}_t, \mathbf{y})} + \tilde{Q}_t^\theta(\mathbf{x}_t, \mathbf{y}) - \tilde{Q}_t^0(\mathbf{x}_t, \mathbf{y}) \right) + o(\Delta t). \quad (19)$$

Summing these contributions and taking the continuous limit ( $\Delta t \rightarrow 0$ ) yields the integral form for the path KL divergence:

$$\mathcal{L}(\theta) = \int_0^1 \mathbb{E}_{\mathbf{x}_t \sim p_{t|0}} \left[ \sum_{\mathbf{y} \neq \mathbf{x}_t} \left( \tilde{Q}_t^0(\mathbf{x}_t, \mathbf{y}) \log \frac{\tilde{Q}_t^0(\mathbf{x}_t, \mathbf{y})}{\tilde{Q}_t^\theta(\mathbf{x}_t, \mathbf{y})} + \tilde{Q}_t^\theta(\mathbf{x}_t, \mathbf{y}) - \tilde{Q}_t^0(\mathbf{x}_t, \mathbf{y}) \right) \right] dt. \quad (20)$$

We now substitute the specific definitions of these transition rates. The true conditional reverse rate  $\tilde{Q}_t^0$  is related to the forward rate  $Q_t(\mathbf{y}, \mathbf{x}_t)$  by

$$\tilde{Q}_t^0(\mathbf{x}_t, \mathbf{y}) = Q_t(\mathbf{y}, \mathbf{x}_t) \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)}. \quad (21)$$

The parameterized reverse rate  $\tilde{Q}_t^\theta$  is defined using the score network  $s_\theta(\mathbf{x}_t, t)_\mathbf{y}$ :

$$\tilde{Q}_t^\theta(\mathbf{x}_t, \mathbf{y}) = Q_t(\mathbf{y}, \mathbf{x}_t) s_\theta(\mathbf{x}_t, t)_\mathbf{y}. \quad (22)$$

We substitute these rates (Eq.21 and Eq.22) into the expression inside the summation in Eq.20. The expression becomes:

$$\left( Q_t(\mathbf{y}, \mathbf{x}_t) \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \right) \log \frac{Q_t(\mathbf{y}, \mathbf{x}_t) \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)}}{Q_t(\mathbf{y}, \mathbf{x}_t) s_\theta(\mathbf{x}_t, t)_\mathbf{y}} + Q_t(\mathbf{y}, \mathbf{x}_t) s_\theta(\mathbf{x}_t, t)_\mathbf{y} - Q_t(\mathbf{y}, \mathbf{x}_t) \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)}. \quad (23)$$

We simplify by canceling  $Q_t(\mathbf{y}, \mathbf{x}_t)$  inside the logarithm and factoring it out from the entire expression:

$$= Q_t(\mathbf{y}, \mathbf{x}_t) \left[ \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \log \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} + s_\theta(\mathbf{x}_t, t)_\mathbf{y} - \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \right]. \quad (24)$$

We expand the logarithm ( $\log(A/B) = \log A - \log B$ ) and rearrange the terms:

$$= Q_t(\mathbf{y}, \mathbf{x}_t) \left[ \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \left( \log \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} - \log s_\theta(\mathbf{x}_t, t)_\mathbf{y} \right) + s_\theta(\mathbf{x}_t, t)_\mathbf{y} - \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \right] \quad (25)$$

$$= Q_t(\mathbf{y}, \mathbf{x}_t) \left[ s_\theta(\mathbf{x}_t, t)_\mathbf{y} - \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \log s_\theta(\mathbf{x}_t, t)_\mathbf{y} + \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \left( \log \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} - 1 \right) \right]. \quad (26)$$

Letting  $K(a) = a(\log a - 1)$ . Recognizing that the time integral  $\int_0^1 dt$  combined with the expectation over  $\mathbf{x}_t \sim p_{t|0}$  is equivalent to the expectation over uniform time  $t \sim U(0, 1)$  and the corresponding conditional  $\mathbf{x}_t$ , the objective  $\mathcal{L}(\theta)$  is exactly the DSE loss (Eq.3):

$$\mathcal{L}_\theta^{\text{DSE}}(\mathbf{x}_0) = \mathbb{E}_{t, \mathbf{x}_t} \sum_{\mathbf{y} \neq \mathbf{x}_t} Q_t(\mathbf{y}, \mathbf{x}_t) \left[ s_\theta(\mathbf{x}_t, t)_\mathbf{y} - \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \log s_\theta(\mathbf{x}_t, t)_\mathbf{y} + K \left( \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \right) \right]. \quad (27)$$

## D.2 DERIVATION OF SEQUENCE-LEVEL TRANSITION PROBABILITY (EQ.5)

We derive the sequence-level transition probability  $p_{t|s}(\mathbf{x}_t|\mathbf{x}_s)$  based on the definition of the DID forward process as an independent token-level deletion process. We begin by analyzing the dynamics of a single token.

The token-level process (Eq.4) is a Continuous-Time Markov Chain (CTMC) on the state space  $\{v, \emptyset\}$ , where  $v \in \mathcal{V}$  denotes the presence of a token and  $\emptyset$  denotes the deleted state. The transition rate matrix  $Q_t^{\text{tok}}$  at time  $t$ , indexed by  $(v, \emptyset)$ , is defined by the deletion rate  $\sigma(t)$ :

$$Q_t^{\text{tok}} = \begin{pmatrix} -\sigma(t) & \sigma(t) \\ 0 & 0 \end{pmatrix}. \quad (28)$$

Let  $P_v(\tau)$  be the probability that the token is in state  $v$  at time  $\tau \in [s, t]$ , given it started in state  $v$  at time  $s$  ( $P_v(s) = 1$ ). The evolution of this probability follows the Kolmogorov forward equation:

$$\frac{dP_v(\tau)}{d\tau} = P_v(\tau)Q_\tau^{\text{tok}}(v, v) + P_\emptyset(\tau)Q_\tau^{\text{tok}}(\emptyset, v) = -\sigma(\tau)P_v(\tau). \quad (29)$$

Solving this first-order ordinary differential equation by integrating from  $s$  to  $t$ :

$$\int_s^t \frac{dP_v(\tau)}{P_v(\tau)} = \int_s^t -\sigma(\tau)d\tau \implies \ln(P_v(t)) - \ln(P_v(s)) = -(\bar{\sigma}(t) - \bar{\sigma}(s)). \quad (30)$$

Thus, the probability that a single token survives during the interval  $[s, t]$  is  $P_v(t) = e^{-(\bar{\sigma}(t) - \bar{\sigma}(s))}$ . Conversely, the probability of deletion is  $1 - e^{-(\bar{\sigma}(t) - \bar{\sigma}(s))}$ .

We now consider the Sequence-level transition from  $\mathbf{x}_s$  to  $\mathbf{x}_t$ . Let  $\Delta\bar{\sigma} = \bar{\sigma}(t) - \bar{\sigma}(s)$ . A transition occurs if the tokens forming  $\mathbf{x}_t$  survive and the remaining  $|\mathbf{x}_s| - |\mathbf{x}_t|$  tokens are deleted. Since token deletions are independent, the probability of a specific path (a specific occurrence of  $\mathbf{x}_t$  in  $\mathbf{x}_s$ ) is  $(e^{-\Delta\bar{\sigma}})^{|\mathbf{x}_t|} \times (1 - e^{-\Delta\bar{\sigma}})^{|\mathbf{x}_s| - |\mathbf{x}_t|}$ .

The total transition probability  $p_{t|s}(\mathbf{x}_t|\mathbf{x}_s)$  is the sum over all distinct paths, counted by the subsequence count  $N(\mathbf{x}_t, \mathbf{x}_s)$ . Therefore, we obtain:

$$p_{t|s}(\mathbf{x}_t|\mathbf{x}_s) = N(\mathbf{x}_t, \mathbf{x}_s)(1 - e^{-(\bar{\sigma}(t) - \bar{\sigma}(s))})^{|\mathbf{x}_s| - |\mathbf{x}_t|} e^{-(\bar{\sigma}(t) - \bar{\sigma}(s))|\mathbf{x}_t|}. \quad (31)$$

### D.3 DERIVATION OF THE TRANSITION RATE (EQ.6)

We derive the transition rate  $Q_t(\mathbf{y}, \mathbf{x}_t)$  from a sequence  $\mathbf{y}$  to a sequence  $\mathbf{x}_t$ , where  $\mathbf{x}_t$  is obtained from  $\mathbf{y}$  by deleting a single token (denoted as  $\mathbf{y} \succ_1 \mathbf{x}_t$ ). This implies  $|\mathbf{y}| = |\mathbf{x}_t| + 1$ .

$$Q_t(\mathbf{y}, \mathbf{x}_t) \triangleq \lim_{\Delta t \rightarrow 0} \frac{p_{t+\Delta t|t}(\mathbf{x}_t|\mathbf{y})}{\Delta t} \quad (32)$$

$$= \lim_{\Delta t \rightarrow 0} \frac{(1 - e^{-(\bar{\sigma}(t+\Delta t) - \bar{\sigma}(t))})^{|\mathbf{y}| - |\mathbf{x}_t|} e^{-(\bar{\sigma}(t+\Delta t) - \bar{\sigma}(t))^{|\mathbf{x}_t|}} N(\mathbf{x}_t, \mathbf{y})}{\Delta t} \quad (33)$$

$$= \lim_{\Delta t \rightarrow 0} \frac{(1 - e^{-\sigma(t)\Delta t + o(\Delta t)})^1 \cdot e^{-(\sigma(t)\Delta t + o(\Delta t))^{|\mathbf{x}_t|}} \cdot N(\mathbf{x}_t, \mathbf{y})}{\Delta t} \quad (34)$$

$$= \lim_{\Delta t \rightarrow 0} \frac{(\sigma(t)\Delta t + o(\Delta t)) \cdot (1 - |\mathbf{x}_t| \sigma(t)\Delta t + o(\Delta t)) \cdot N(\mathbf{x}_t, \mathbf{y})}{\Delta t} \quad (35)$$

$$= \lim_{\Delta t \rightarrow 0} \frac{\sigma(t)\Delta t + o(\Delta t)}{\Delta t} \cdot N(\mathbf{x}_t, \mathbf{y}) \quad (36)$$

$$= \sigma(t)N(\mathbf{x}_t, \mathbf{y}). \quad (37)$$

### D.4 DERIVATION OF THE RELATIONSHIP BETWEEN CONCRETE SCORE AND INSERTION SCORE (EQ. 8)

We aim to prove the identity stated in Eq. 8. This identity concerns the backward process where transitions occur such that  $\mathbf{y} \succ_1 \mathbf{x}_t$ . Note that this condition implies  $|\mathbf{y}| = |\mathbf{x}_t| + 1$ .

First, we derive the explicit form of the concrete score  $s(\mathbf{x}_t, t)_{\mathbf{y}} = p_t(\mathbf{y})/p_t(\mathbf{x}_t)$  by expanding the marginal distributions using the forward transition probability (Eq. 5):

$$\begin{aligned} s(\mathbf{x}_t, t)_{\mathbf{y}} &= \frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)} = \frac{\mathbb{E}_{\mathbf{x}_0}[p_{t|0}(\mathbf{y}|\mathbf{x}_0)]}{\mathbb{E}_{\mathbf{x}_0}[p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)]} \\ &= \frac{\mathbb{E}_{\mathbf{x}_0}[(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0| - |\mathbf{y}|} e^{-\bar{\sigma}(t)|\mathbf{y}|} N(\mathbf{y}, \mathbf{x}_0)]}{\mathbb{E}_{\mathbf{x}_0}[(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0| - |\mathbf{x}_t|} e^{-\bar{\sigma}(t)|\mathbf{x}_t|} N(\mathbf{x}_t, \mathbf{x}_0)]} \\ &= \frac{e^{-\bar{\sigma}(t)|\mathbf{y}|} (1 - e^{-\bar{\sigma}(t)})^{-|\mathbf{y}|} \mathbb{E}_{\mathbf{x}_0}[(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0|} N(\mathbf{y}, \mathbf{x}_0)]}{e^{-\bar{\sigma}(t)|\mathbf{x}_t|} (1 - e^{-\bar{\sigma}(t)})^{-|\mathbf{x}_t|} \mathbb{E}_{\mathbf{x}_0}[(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0|} N(\mathbf{x}_t, \mathbf{x}_0)]} \\ &\stackrel{|\mathbf{y}| = |\mathbf{x}_t| + 1}{=} \frac{e^{-\bar{\sigma}(t)} \mathbb{E}_{\mathbf{x}_0}[(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0|} N(\mathbf{y}, \mathbf{x}_0)]}{1 - e^{-\bar{\sigma}(t)} \mathbb{E}_{\mathbf{x}_0}[(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0|} N(\mathbf{x}_t, \mathbf{x}_0)]}. \end{aligned} \quad (38)$$

Next, we examine the right-hand side (RHS) of Eq. 8 and substitute the definition of the insertion score  $\bar{s}$  (Eq. 7):

$$\text{RHS} = \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \frac{1}{N(\mathbf{x}_t, \mathbf{y})} \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} \bar{s}(\mathbf{x}_t, t)[i, v(\mathbf{x}_t, \mathbf{y})] \quad (39)$$

$$= \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \frac{1}{N(\mathbf{x}_t, \mathbf{y})} \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} \left( \frac{\mathbb{E}_{\mathbf{x}_0}[(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0|} N(\text{Ins}(\mathbf{x}_t, i, v(\mathbf{x}_t, \mathbf{y})), \mathbf{x}_0)]}{\mathbb{E}_{\mathbf{x}_0}[(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0|} N(\mathbf{x}_t, \mathbf{x}_0)]} \right). \quad (40)$$

By definition of the index set  $I(\mathbf{x}_t, \mathbf{y})$ , for any  $i \in I(\mathbf{x}_t, \mathbf{y})$ , we have  $\text{Ins}(\mathbf{x}_t, i, v(\mathbf{x}_t, \mathbf{y})) = \mathbf{y}$ . Therefore:

$$\text{RHS} = \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \frac{1}{N(\mathbf{x}_t, \mathbf{y})} \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} \left( \frac{\mathbb{E}_{\mathbf{x}_0}[(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0|} N(\mathbf{y}, \mathbf{x}_0)]}{\mathbb{E}_{\mathbf{x}_0}[(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0|} N(\mathbf{x}_t, \mathbf{x}_0)]} \right). \quad (41)$$

The term inside the summation is independent of the index  $i$ . We factor it out and utilize the property that  $\sum_{i \in I(\mathbf{x}_t, \mathbf{y})} 1 = |I(\mathbf{x}_t, \mathbf{y})| = N(\mathbf{x}_t, \mathbf{y})$ :

$$\text{RHS} = \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \frac{1}{N(\mathbf{x}_t, \mathbf{y})} \left( \frac{\mathbb{E}_{\mathbf{x}_0} [(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0|} N(\mathbf{y}, \mathbf{x}_0)]}{\mathbb{E}_{\mathbf{x}_0} [(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0|} N(\mathbf{x}_t, \mathbf{x}_0)]} \right) N(\mathbf{x}_t, \mathbf{y}) \quad (42)$$

$$= \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \frac{\mathbb{E}_{\mathbf{x}_0} [(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0|} N(\mathbf{y}, \mathbf{x}_0)]}{\mathbb{E}_{\mathbf{x}_0} [(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0|} N(\mathbf{x}_t, \mathbf{x}_0)]}. \quad (43)$$

This matches the derived concrete score in Eq. 38, completing the proof.

## D.5 DERIVATION OF THE SAMPLING PROBABILITY (EQ.10)

We aim to derive the probability of executing a specific insertion operation—inserting token  $v$  after position  $i$ —within an infinitesimal time interval  $[t - \Delta t, t]$  during the backward process.

The backward process is a CTMC characterized by the parameterized reverse transition rate matrix  $\tilde{Q}_t^\theta$ . The rate of transition from state  $\mathbf{x}_t$  to a different state  $\mathbf{y}$  is defined as:

$$\tilde{Q}_t^\theta(\mathbf{x}_t, \mathbf{y}) = Q_t(\mathbf{y}, \mathbf{x}_t) s_\theta(\mathbf{x}_t, t)_{\mathbf{y}}. \quad (44)$$

In the DID framework, backward transitions occur only when  $\mathbf{y} \succ_1 \mathbf{x}_t$ . We substitute the forward transition rate (Eq.6),  $Q_t(\mathbf{y}, \mathbf{x}_t) = \sigma(t)N(\mathbf{x}_t, \mathbf{y})$ . We also substitute the parameterized concrete score  $s_\theta$ , which is derived by parameterizing the relationship between the concrete score and the insertion score (Eq.8):

$$s_\theta(\mathbf{x}_t, t)_{\mathbf{y}} = \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \frac{1}{N(\mathbf{x}_t, \mathbf{y})} \sum_{j \in I(\mathbf{x}_t, \mathbf{y})} \bar{s}_\theta(\mathbf{x}_t, t)[j, v(\mathbf{x}_t, \mathbf{y})]. \quad (45)$$

Substituting these expressions into the definition of  $\tilde{Q}_t^\theta(\mathbf{x}_t, \mathbf{y})$ :

$$\tilde{Q}_t^\theta(\mathbf{x}_t, \mathbf{y}) = (\sigma(t)N(\mathbf{x}_t, \mathbf{y})) \cdot \left( \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \frac{1}{N(\mathbf{x}_t, \mathbf{y})} \sum_{j \in I(\mathbf{x}_t, \mathbf{y})} \bar{s}_\theta(\mathbf{x}_t, t)[j, v(\mathbf{x}_t, \mathbf{y})] \right) \quad (46)$$

$$= \sigma(t) \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{j \in I(\mathbf{x}_t, \mathbf{y})} \bar{s}_\theta(\mathbf{x}_t, t)[j, v(\mathbf{x}_t, \mathbf{y})]. \quad (47)$$

This result demonstrates that the total transition rate from  $\mathbf{x}_t$  to  $\mathbf{y}$  is decomposed into a summation of individual components. Each term in the summation,  $\sigma(t) \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \bar{s}_\theta(\mathbf{x}_t, t)[j, v(\mathbf{x}_t, \mathbf{y})]$ , corresponds to the instantaneous rate of the specific insertion operation  $(j, v(\mathbf{x}_t, \mathbf{y}))$  that transforms  $\mathbf{x}_t$  into  $\mathbf{y}$ .

By the definition of a CTMC, the probability of a specific operation  $(i, v)$  occurring within the infinitesimal interval  $\Delta t$  is given by its corresponding rate multiplied by  $\Delta t$ . For  $v \neq \emptyset$ , we identify this probability directly from the decomposition above:

$$p_{t-\Delta t|t}^\theta((i, v)|\mathbf{x}_t) = \left( \sigma(t) \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \bar{s}_\theta(\mathbf{x}_t, t)[i, v] \right) \Delta t + o(\Delta t) \quad (48)$$

$$= \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \bar{s}_\theta(\mathbf{x}_t, t)[i, v] \Delta t + o(\Delta t). \quad (49)$$

This confirms the first case of Eq.10. The probability of no insertion occurring at position  $i$  (i.e.,  $v = \emptyset$ ) is determined by the normalization constraint, ensuring the sum of probabilities for all possible events at that position equals 1.

To confirm the equivalence between this action-based sampling and the required state-based transition, we verify that the action probabilities correctly recover the state transition probability  $p_{t-\Delta t|t}^\theta(\mathbf{y}|\mathbf{x}_t)$ . A transition to  $\mathbf{y}$  occurs if any of the actions indexed by  $I(\mathbf{x}_t, \mathbf{y})$  occurs. In the infinitesimal limit  $\Delta t \rightarrow 0$ , these actions are mutually exclusive events:

$$p_{t-\Delta t|t}^\theta(\mathbf{y}|\mathbf{x}_t) = \sum_{j \in I(\mathbf{x}_t, \mathbf{y})} p_{t-\Delta t|t}^\theta((j, v(\mathbf{x}_t, \mathbf{y}))|\mathbf{x}_t) + o(\Delta t) \quad (50)$$

$$= \left( \sigma(t) \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{j \in I(\mathbf{x}_t, \mathbf{y})} \bar{s}_\theta(\mathbf{x}_t, t)[j, v(\mathbf{x}_t, \mathbf{y})] \right) \Delta t + o(\Delta t). \quad (51)$$

By Eq. 47, the term in the parenthesis is exactly  $\tilde{Q}_t^\theta(\mathbf{x}_t, \mathbf{y})$ . Thus,  $p_{t-\Delta t|t}^\theta(\mathbf{y}|\mathbf{x}_t) = \tilde{Q}_t^\theta(\mathbf{x}_t, \mathbf{y})\Delta t + o(\Delta t)$ , confirming that the action-based sampling correctly implements the dynamics of the backward CTMC.

#### D.6 DERIVATION OF THE DISE OBJECTIVE (EQ.11)

We derive the Denoising Insertion Score Entropy (DISE) objective (Eq.11) starting from the general DSE objective (Eq.3), demonstrating that DISE is a variational upper bound on DSE,  $\mathcal{L}_\theta^{\text{DISE}}(\mathbf{x}_0) \geq \mathcal{L}_\theta^{\text{DSE}}(\mathbf{x}_0)$ .

We begin with the DSE objective:

$$\mathcal{L}_\theta^{\text{DSE}}(\mathbf{x}_0) = \mathbb{E}_{t, \mathbf{x}_t} \sum_{\mathbf{y} \neq \mathbf{x}_t} Q_t(\mathbf{y}, \mathbf{x}_t) \left[ s_\theta(\mathbf{x}_t, t)_{\mathbf{y}} - \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \log s_\theta(\mathbf{x}_t, t)_{\mathbf{y}} + K \left( \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \right) \right], \quad (52)$$

where  $K(a) = a(\log a - 1)$ . In the deletion–insertion process, the transition rate  $Q_t(\mathbf{y}, \mathbf{x}_t)$  is non-zero only when  $\mathbf{y} \succ_1 \mathbf{x}_t$ .

We now substitute the specific definitions for the DID process. The transition rate is  $Q_t(\mathbf{y}, \mathbf{x}_t) = \sigma(t)N(\mathbf{x}_t, \mathbf{y})$ . The conditional probability ratio is:

$$\frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} = \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \frac{N(\mathbf{y}, \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)}. \quad (53)$$

The parameterized concrete score (from parameterized Eq.8) is:

$$s_\theta(\mathbf{x}_t, t)_{\mathbf{y}} = \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \frac{1}{N(\mathbf{x}_t, \mathbf{y})} \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} \bar{s}_\theta(\mathbf{x}_t, t)[i, v(\mathbf{x}_t, \mathbf{y})]. \quad (54)$$

We examine the expression inside the bracket of the DSE objective by substituting these definitions. Let  $v = v(\mathbf{x}_t, \mathbf{y})$  for brevity in the following block:

$$\begin{aligned} & s_\theta(\mathbf{x}_t, t)_{\mathbf{y}} - \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \log s_\theta(\mathbf{x}_t, t)_{\mathbf{y}} + K \left( \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \right) \\ &= \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \frac{1}{N(\mathbf{x}_t, \mathbf{y})} \left( \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} \bar{s}_\theta(\mathbf{x}_t, t)[i, v] \right) \\ & \quad - \left( \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \frac{N(\mathbf{y}, \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \right) \log \left( \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \frac{1}{N(\mathbf{x}_t, \mathbf{y})} \left( \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} \bar{s}_\theta(\mathbf{x}_t, t)[i, v] \right) \right) \\ & \quad + K \left( \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \frac{N(\mathbf{y}, \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \right). \end{aligned} \quad (55)$$

We expand the  $K(\cdot)$  term using  $K(a) = a(\log a - 1)$ . By expanding the logarithms ( $\log(AB) = \log A + \log B$ ), we observe that the terms involving the time factor  $\log\left(\frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}}\right)$  cancel out exactly.

The expression inside the bracket simplifies significantly by factoring out the time prefactor:

$$\begin{aligned}
&= \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \left[ \frac{1}{N(\mathbf{x}_t, \mathbf{y})} \left( \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} \bar{s}_\theta(\mathbf{x}_t, t)[i, v] \right) \right. \\
&\quad \left. - \frac{N(\mathbf{y}, \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \log \left( \frac{1}{N(\mathbf{x}_t, \mathbf{y})} \left( \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} \bar{s}_\theta(\mathbf{x}_t, t)[i, v] \right) \right) \right] \\
&\quad + \frac{N(\mathbf{y}, \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \left( \log \frac{N(\mathbf{y}, \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} - 1 \right). \tag{56}
\end{aligned}$$

Substituting this simplified bracket back into the main objective equation and multiplying by  $Q_t(\mathbf{y}, \mathbf{x}_t) = \sigma(t)N(\mathbf{x}_t, \mathbf{y})$ . The DSE objective can be decomposed into three terms (T1, T2, T3):

$$\mathcal{L}_\theta^{\text{DSE}}(\mathbf{x}_0) = \text{T1} + \text{T2} + \text{T3}. \tag{57}$$

We define these terms based on the components derived above:

$$\text{T1} = \mathbb{E}_{t, \mathbf{x}_t} \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{\mathbf{y} \succ_1 \mathbf{x}_t} N(\mathbf{x}_t, \mathbf{y}) \left[ \frac{1}{N(\mathbf{x}_t, \mathbf{y})} \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} \bar{s}_\theta(\mathbf{x}_t, t)[i, v] \right]. \tag{58}$$

$$\text{T2} = \mathbb{E}_{t, \mathbf{x}_t} \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{\mathbf{y} \succ_1 \mathbf{x}_t} N(\mathbf{x}_t, \mathbf{y}) \left[ -\frac{N(\mathbf{y}, \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \log \left( \frac{1}{N(\mathbf{x}_t, \mathbf{y})} \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} \bar{s}_\theta(\mathbf{x}_t, t)[i, v] \right) \right]. \tag{59}$$

$$\text{T3} = \mathbb{E}_{t, \mathbf{x}_t} \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{\mathbf{y} \succ_1 \mathbf{x}_t} N(\mathbf{x}_t, \mathbf{y}) K \left( \frac{N(\mathbf{y}, \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \right). \tag{60}$$

We now apply Jensen’s inequality to T2. The term inside the logarithm is an average of insertion scores. Because the logarithm function is concave,  $\log\left(\frac{1}{N} \sum a_i\right) \geq \frac{1}{N} \sum \log a_i$ . Since the logarithm is negated, this leads to an upper bound on T2:

$$-\log \left( \frac{1}{N(\mathbf{x}_t, \mathbf{y})} \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} \bar{s}_\theta(\mathbf{x}_t, t)[i, v] \right) \leq -\frac{1}{N(\mathbf{x}_t, \mathbf{y})} \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} \log \bar{s}_\theta(\mathbf{x}_t, t)[i, v]. \tag{61}$$

We define  $\text{T2}_{\text{Bound}}$  as the upper bound for T2:

$$\text{T2} \leq \text{T2}_{\text{Bound}} = \mathbb{E}_{t, \mathbf{x}_t} \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{\mathbf{y} \succ_1 \mathbf{x}_t} N(\mathbf{x}_t, \mathbf{y}) \frac{N(\mathbf{y}, \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \left[ -\frac{1}{N(\mathbf{x}_t, \mathbf{y})} \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} \log \bar{s}_\theta(\mathbf{x}_t, t)[i, v] \right]. \tag{62}$$

We define the DISE objective as the upper bound obtained by replacing T2 with  $\text{T2}_{\text{Bound}}$ , ensuring  $\mathcal{L}_\theta^{\text{DISE}}(\mathbf{x}_0) \geq \mathcal{L}_\theta^{\text{DSE}}(\mathbf{x}_0)$ :

$$\mathcal{L}_\theta^{\text{DISE}}(\mathbf{x}_0) = \text{T1} + \text{T2}_{\text{Bound}} + \text{T3}. \tag{63}$$

To simplify the nested summations and transform the objective from state-level ( $\mathbf{y}$ ) to operation-level ( $(i, v)$ ), we rely on the following identity.

**Lemma 1** (Summation change of variables). Let  $v(\mathbf{x}_t, \mathbf{y})$  denote the unique inserted token and  $I(\mathbf{x}_t, \mathbf{y})$  the set of valid insertion positions that yield  $\mathbf{y}$  from  $\mathbf{x}_t$ . For any function  $G$ ,

$$\sum_{\mathbf{y} \succ_1 \mathbf{x}_t} \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} G(i, v(\mathbf{x}_t, \mathbf{y}), \mathbf{y}) = \sum_{i, v} G(i, v, \text{Ins}(\mathbf{x}_t, i, v)). \quad (64)$$

*Proof.* Define the index sets

$$\mathcal{A} = \{(\mathbf{y}, i) : \mathbf{y} \succ_1 \mathbf{x}_t, i \in I(\mathbf{x}_t, \mathbf{y})\}, \quad \mathcal{B} = \{(i, v) : i \in \{0, \dots, |\mathbf{x}_t|\}, v \in \mathcal{V}\}. \quad (65)$$

The map  $g : \mathcal{B} \rightarrow \mathcal{A}$  defined by  $g(i, v) = (\text{Ins}(\mathbf{x}_t, i, v), i)$  is a bijection, with its inverse  $f : \mathcal{A} \rightarrow \mathcal{B}$  defined by  $f(\mathbf{y}, i) = (i, v(\mathbf{x}_t, \mathbf{y}))$ . Changing variables over this bijection gives the claimed identity.  $\square$

We apply Lemma 1 to simplify T1, T2<sub>Bound</sub>, and T3.

For T1, the  $N(\mathbf{x}_t, \mathbf{y})$  terms cancel before the summation.

$$\text{T1} = \mathbb{E}_{t, \mathbf{x}_t} \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{\mathbf{y} \succ_1 \mathbf{x}_t} \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} \bar{s}_\theta(\mathbf{x}_t, t)[i, v(\mathbf{x}_t, \mathbf{y})] \quad (66)$$

$$= \mathbb{E}_{t, \mathbf{x}_t} \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{i, v} \bar{s}_\theta(\mathbf{x}_t, t)[i, v]. \quad (\text{Using Lemma 1}) \quad (67)$$

For T2<sub>Bound</sub>, cancellation of  $N(\mathbf{x}_t, \mathbf{y})$  yields:

$$\text{T2}_{\text{Bound}} = \mathbb{E}_{t, \mathbf{x}_t} \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{\mathbf{y} \succ_1 \mathbf{x}_t} \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} -\frac{N(\mathbf{y}, \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \log \bar{s}_\theta(\mathbf{x}_t, t)[i, v(\mathbf{x}_t, \mathbf{y})]. \quad (68)$$

Using Lemma 1, noting that  $\mathbf{y} = \text{Ins}(\mathbf{x}_t, i, v)$ :

$$\text{T2}_{\text{Bound}} = \mathbb{E}_{t, \mathbf{x}_t} \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{i, v} -\frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \log \bar{s}_\theta(\mathbf{x}_t, t)[i, v]. \quad (69)$$

For T3, we first utilize the fact that  $N(\mathbf{x}_t, \mathbf{y}) = \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} 1$ .

$$\text{T3} = \mathbb{E}_{t, \mathbf{x}_t} \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{\mathbf{y} \succ_1 \mathbf{x}_t} \sum_{i \in I(\mathbf{x}_t, \mathbf{y})} K \left( \frac{N(\mathbf{y}, \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \right). \quad (70)$$

Applying Lemma 1 to T3:

$$\text{T3} = \mathbb{E}_{t, \mathbf{x}_t} \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{i, v} K \left( \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \right). \quad (71)$$

We combine T1, T2<sub>Bound</sub>, and T3 to obtain the final DISE objective. Let  $C$  denote the constant term T3's summand.

$$\mathcal{L}_\theta^{\text{DISE}}(\mathbf{x}_0) = \text{T1} + \text{T2}_{\text{Bound}} + \text{T3} \quad (72)$$

$$= \mathbb{E}_{t, \mathbf{x}_t} \left\{ \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{i, v} \left[ \bar{s}_\theta(\mathbf{x}_t, t)[i, v] - \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \log \bar{s}_\theta(\mathbf{x}_t, t)[i, v] + C \right] \right\}. \quad (73)$$

## D.7 DERIVATION OF THE SEQUENCE-LEVEL NORMALIZATION PROPERTY (EQ. 16)

In the fixed-length setting (Sec. 3.5), we assume  $|\mathbf{x}_0| = K$  (a constant) for all data samples. Under this assumption, the insertion score becomes time-independent. We aim to prove the normalization property (Eq. 16), restated here for the fixed-length context:

$$\sum_{i,v} \bar{s}(\mathbf{x}_t)[i, v] = K - |\mathbf{x}_t|. \quad (74)$$

The proof relies on a fundamental combinatorial identity regarding subsequence counts.

**Lemma 2** (Subsequence Count Identity). For any sequences  $\mathbf{x}_t$  and  $\mathbf{x}_0$  such that  $\mathbf{x}_t$  is a subsequence of  $\mathbf{x}_0$ , the following identity holds:

$$\sum_{i,v} N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0) = N(\mathbf{x}_t, \mathbf{x}_0)(|\mathbf{x}_0| - |\mathbf{x}_t|). \quad (75)$$

*Proof of Lemma 2.* We prove the statement  $\sum_{i,v} N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0) = N(\mathbf{x}_t, \mathbf{x}_0)(|\mathbf{x}_0| - |\mathbf{x}_t|)$  via a bijective proof, by constructing two sets of equal cardinality. Let  $S(\mathbf{x}, \mathbf{z})$  denote the set of index tuples corresponding to all occurrences of a subsequence  $\mathbf{x}$  in  $\mathbf{z}$ , such that  $N(\mathbf{x}, \mathbf{z}) = |S(\mathbf{x}, \mathbf{z})|$ .

First, consider the set A, defined as the set of pairs  $(I, j)$ , where  $I$  is the index tuple of an occurrence of  $\mathbf{x}_t$  in  $\mathbf{x}_0$ , and  $j$  is an index in  $\mathbf{x}_0$  that is not part of that occurrence:

$$A = \{(I, j) : I \in S(\mathbf{x}_t, \mathbf{x}_0), j \in \{1, \dots, |\mathbf{x}_0|\} \setminus I\}. \quad (76)$$

The cardinality of A is  $|A| = N(\mathbf{x}_t, \mathbf{x}_0)(|\mathbf{x}_0| - |\mathbf{x}_t|)$ .

Second, consider the set B, defined as the set of pairs  $((i, v), J)$ , where  $(i, v)$  is an insertion operation on  $\mathbf{x}_t$ , and  $J$  is the index tuple of an occurrence of the resulting sequence,  $\text{Ins}(\mathbf{x}_t, i, v)$ , in  $\mathbf{x}_0$ :

$$B = \{((i, v), J) : J \in S(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)\}. \quad (77)$$

The cardinality of B is  $|B| = \sum_{i,v} N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)$ .

We now establish a bijection between A and B. For any element  $(I, j) \in A$ , we define a mapping to an element in B as follows: let the inserted token be  $v = \mathbf{x}_0[j]$ , and let the insertion position relative to  $\mathbf{x}_t$  be  $i = |\{k \in I : k < j\}|$ . The new index tuple is  $J = I \cup \{j\}$ , sorted. The subsequence  $\mathbf{x}_0[J]$  is precisely  $\text{Ins}(\mathbf{x}_t, i, v)$  by construction. This defines a unique mapping  $f : A \rightarrow B$ .

Conversely, for any element  $((i, v), J) \in B$ , we can define an inverse mapping. The index of the inserted token in  $\mathbf{x}_0$  is the  $(i + 1)$ -th element of the sorted tuple  $J$ ; let this be  $j$ . Removing this index yields the tuple  $I = J \setminus \{j\}$ , which corresponds to an occurrence of  $\mathbf{x}_t$ . This defines a unique mapping  $g : B \rightarrow A$ .

Since a one-to-one correspondence exists between the sets, their cardinalities must be equal. Therefore,  $|A| = |B|$ , which proves the lemma.  $\square$

*Proof of Eq. 16.* Under the fixed-length assumption ( $|\mathbf{x}_0| = K$ ), the definition of the insertion score (Eq. 7) simplifies because the time-dependent terms  $(1 - e^{-\bar{\sigma}(t)})^{|\mathbf{x}_0|}$  are constant  $(1 - e^{-\bar{\sigma}(t)})^K$  and cancel out, leading to the time-independent score:

$$\bar{s}(\mathbf{x}_t)[i, v] = \frac{\mathbb{E}_{\mathbf{x}_0}[N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)]}{\mathbb{E}_{\mathbf{x}_0}[N(\mathbf{x}_t, \mathbf{x}_0)]}. \quad (78)$$

We sum this score over all possible insertion operations  $(i, v)$ :

$$\sum_{i,v} \bar{s}(\mathbf{x}_t)[i, v] = \sum_{i,v} \frac{\mathbb{E}_{\mathbf{x}_0}[N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)]}{\mathbb{E}_{\mathbf{x}_0}[N(\mathbf{x}_t, \mathbf{x}_0)]} \quad (79)$$



$$= \frac{1}{\mathbb{E}_{\mathbf{x}_0}[N(\mathbf{x}_t, \mathbf{x}_0)]} \mathbb{E}_{\mathbf{x}_0} \left[ \sum_{i,v} N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0) \right]. \quad (80)$$

We apply the combinatorial identity (Lemma 2) to the summation inside the expectation, substituting  $|\mathbf{x}_0| = K$ :

$$\sum_{i,v} N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0) = N(\mathbf{x}_t, \mathbf{x}_0)(K - |\mathbf{x}_t|). \quad (81)$$

Substituting this back:

$$\sum_{i,v} \bar{s}(\mathbf{x}_t)[i, v] = \frac{1}{\mathbb{E}_{\mathbf{x}_0}[N(\mathbf{x}_t, \mathbf{x}_0)]} \mathbb{E}_{\mathbf{x}_0} [N(\mathbf{x}_t, \mathbf{x}_0)(K - |\mathbf{x}_t|)]. \quad (82)$$

Since  $(K - |\mathbf{x}_t|)$  is constant with respect to the expectation over  $\mathbf{x}_0$ :

$$\sum_{i,v} \bar{s}(\mathbf{x}_t)[i, v] = (K - |\mathbf{x}_t|) \frac{\mathbb{E}_{\mathbf{x}_0}[N(\mathbf{x}_t, \mathbf{x}_0)]}{\mathbb{E}_{\mathbf{x}_0}[N(\mathbf{x}_t, \mathbf{x}_0)]} = K - |\mathbf{x}_t|. \quad (83)$$

This confirms the normalization property stated in Eq. 16.  $\square$

#### D.8 DERIVATION OF THE DICE OBJECTIVE FOR FIXED-LENGTH DATA (EQ. 17)

In the fixed-length setting (Section 3.5), we assume  $|\mathbf{x}_0| = K$  (constant). This assumption leads to time-independent insertion scores  $\bar{s}_\theta(\mathbf{x}_t)[i, v]$  (as shown in Appendix D.7) and allows for the exact simplification of the DISE objective into the Denoising Insertion Cross-Entropy (DICE) objective.

We start from the DISE objective (Eq.11):

$$\begin{aligned} \mathcal{L}_\theta^{\text{DISE}}(\mathbf{x}_0) = \mathbb{E}_{t, \mathbf{x}_t} \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{i,v} \left[ \bar{s}_\theta(\mathbf{x}_t)[i, v] - \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \log \bar{s}_\theta(\mathbf{x}_t)[i, v] \right. \\ \left. + K \left( \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \right) \right]. \quad (84) \end{aligned}$$

We rearrange the expression inside the square brackets using the definition  $K(a) = a(\log a - 1)$ .

$$\begin{aligned} \bar{s}_\theta[i, v] - \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \log \bar{s}_\theta[i, v] + \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \left( \log \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} - 1 \right) \\ = \left( \bar{s}_\theta[i, v] - \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \right) \\ + \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \left( \log \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} - \log \bar{s}_\theta[i, v] \right). \quad (85) \end{aligned}$$

We substitute this rearrangement back into the DISE objective:

$$\begin{aligned} \mathcal{L}_\theta^{\text{DISE}}(\mathbf{x}_0) = \mathbb{E}_{t, \mathbf{x}_t} \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{i,v} \left[ \left( \bar{s}_\theta[i, v] - \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \right) \right. \\ \left. + \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \left( \log \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} - \log \bar{s}_\theta[i, v] \right) \right]. \quad (86) \end{aligned}$$

We now utilize the crucial normalization properties derived from the fixed-length assumption ( $|\mathbf{x}_0| = K$ ). From Lemma 2 (Appendix D.7), the true subsequence count ratios satisfy:

$$\sum_{i,v} \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} = K - |\mathbf{x}_t|. \quad (87)$$

As discussed in Section 3.5, we design the network architecture such that the parameterized scores  $\bar{s}_\theta$  exactly satisfy the same normalization constraint:

$$\sum_{i,v} \bar{s}_\theta(\mathbf{x}_t)[i, v] = K - |\mathbf{x}_t|. \quad (88)$$

Because both the true ratios and the parameterized scores sum to the same value, the summation of the first group of terms in Eq. 85 vanishes:

$$\sum_{i,v} \left( \bar{s}_\theta(\mathbf{x}_t)[i, v] - \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \right) = (K - |\mathbf{x}_t|) - (K - |\mathbf{x}_t|) = 0. \quad (89)$$

Therefore, the DISE objective simplifies exactly. We define this simplified form as the DICE objective, which is exactly equal to the DISE loss under the fixed-length setting ( $\mathcal{L}_\theta^{\text{DICE}}(\mathbf{x}_0) = \mathcal{L}_\theta^{\text{DISE}}(\mathbf{x}_0)$ ):

$$\mathcal{L}_\theta^{\text{DICE}}(\mathbf{x}_0) = \mathbb{E}_{t, \mathbf{x}_t} \frac{\sigma(t) e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \sum_{i,v} \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \left( \log \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} - \log \bar{s}_\theta(\mathbf{x}_t)[i, v] \right). \quad (90)$$

Rearranging the terms inside the summation yields the final DICE objective:

$$\mathcal{L}_\theta^{\text{DICE}}(\mathbf{x}_0) = \mathbb{E}_{t, \mathbf{x}_t} \left\{ \sum_{i,v} \frac{\sigma(t) e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \left[ -\log \bar{s}_\theta(\mathbf{x}_t)[i, v] + C \right] \right\}, \quad (91)$$

where  $C = \log \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)}$  is a  $\theta$ -free constant.

## E EXPERIMENTAL DETAILS

### E.1 FIXED-LENGTH TRAINING DETAILS

**Model architecture.** Following RADD, our models use an encoder-only transformer architecture with a dropout rate 0.02, rotary position embedding (Su et al., 2023), and untied word embeddings between input and output. Other model architecture hyperparameters are summarized in Tab. 7. The FFN dimension is  $4 \times$  hidden dimension.

Size	Layers	Hidden Dimension	Attention Heads	Non-Embedding Parameters
Small	12	768	12	85M
Medium	24	1024	16	302M
Large	36	1280	20	708M

Table 7: Model architecture hyperparameters for difference model sizes of RADD and DID.

We use FlashAttention (Dao et al., 2022) to support attention computation for packed variable-length sequences, while RADD uses PyTorch’s standard scaled dot product attention for its fixed-length batched inputs, which is also based on FlashAttention.

**Dataset and pre-processing.** For a fair comparison, we train DID on OpenWebText dataset (Gokaslan & Cohen, 2019) with the same steps (400K steps) or compute budget (800K steps) as RADD (Ou et al., 2025), under the log-linear noise schedule  $\bar{\sigma}(t) = -\log(1 - t)$ , we train DID with a batch size 512, and a context length 1024. We tokenize OpenWebText dataset with the GPT2 (Radford et al., 2019) tokenizer, the vocabulary size is 50257. We follow the data pre-processing adopted in RADD, concatenating all sequences in the training dataset, then splitting it into chunks of fixed length 1024. We add a `<BOS>` special token to the first position of each chunk to model the insertion behavior before the first normal token by modeling the insertion after the `<BOS>` special token. We use the same token as `<EOS>` for `<BOS>`. We evaluate the zero-shot language modeling perplexity on WikiText, Lambada, Scientific Papers (Arxiv and Pubmed, abstract parts), AG News, LM1B, and PTB datasets (Merity et al., 2017; Paperno et al., 2016; Cohan et al., 2018; Zhang et al., 2015; Chelba et al., 2014; Marcus et al., 1993) in Tab. 1.

**Optimization.** For a fair comparison, we did not do a hyperparameter search, our optimization configuration strictly follows RADD, we use AdamW (Loshchilov & Hutter, 2019) optimizer with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 1e-8$ , a weight decay rate of 0.03, a constant learning rate  $3e-4$  that linearly warmed up from 0 over the first 2500 steps, a gradient norm clipping value 1, an exponential moving average (EMA) with a decay rate 0.9999, and float16 is enabled for mixed-precision training. To avoid OOM errors during model training, we set the gradient accumulation steps to 1, 2, 2 for DID small, medium, and large; and 2, 4, 4 for RADD small, medium, and large.

## E.2 VARIABLE-LENGTH TRAINING DETAILS

**Model architecture.** Following ILM (Patel et al., 2025), we train a small model with 85M non-embedding parameters. Different from the fixed-length model, the variable-length model is not time-independent, i.e. we need to input the time information into the network. Therefore, we employ an adaptive layernorm for each transformer block as in the practice of DiT (Peebles & Xie, 2023) (as well as the time-dependent masked diffusion models like SEDD (Lou et al., 2024) and MDLM (Sahoo et al., 2024)), the condition embedding dimension is 128, which brings about extra parameters. To keep the total parameter amount of a transformer block unchanged, we reduce the FFN dimension to  $3.5 \times$  hidden dimension. Following ILM, we use a dropout rate of 0.1, rotary positional embedding, and untied word embeddings between input and output. FlashAttention is also employed for the efficient computation of variable-length data.

**Dataset and pre-processing.** Following ILM, we train variable-length models on the Stories dataset. We tokenize it with the Bert-Base-Uncased tokenizer, whose vocabulary size is 30522. Each sentence in the datasets is an individual training datapoint; hence, the training sequences are of variable lengths. The Stories dataset is truncated with a maximum context length of 1024. The batch size is 512. Data is also noised with a log-linear noise schedule in the diffusion forward process.

**Optimization.** Following the experiment settings of ILM, we use an AdamW optimizer with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 1e-8$ , a weight decay of 0.01 on all parameters (including biases and normalization layers), a constant learning rate  $1e-4$  that linearly warmed up from 0 over the first 1000 steps, a gradient norm clipping value 1, an exponential moving average (EMA) with a decay rate 0.9999, and bfloat16 is enabled for mixed-precision training. Stories dataset is trained for 60k steps. To avoid OOM errors during model training, we set the gradient accumulation steps to 1, 1, 1 for DID small, medium, and large; and 2, 4, 4 for RADD small, medium, and large.

### 1269 E.3 SAMPLING DETAILS

1270  
1271 **Timestep discretization.** We use a standard uniform timestep discretization grid, i.e.  $t(i) = i/N$ , where  $N$   
1272 is the total number of denoising steps in generation.

1273 **The number of samples.** We evaluate the averaged generative perplexity (evaluated by GPT2 Large), unigram  
1274 entropy, generation length, and inference time over 1024 samples with a batch size of 32.

1275  
1276 **Data precision.** According to the precision issue of Gumbel-argmax sampling with float32 discussed in (Ou  
1277 et al., 2025; Zheng et al., 2025), we use float64 precision for all generation tasks to ensure accurate categorical  
1278 samplings.

### 1279 E.4 EVALUATION METRICS

1280  
1281 **Zero-shot language modeling perplexity** is used to evaluate how well a model is trained, which uses the  
1282 model to be evaluated to calculate the exponential of the average negative log-likelihood per token on datasets  
1283 unseen by the model (i.e. zero-shot):

$$1284 \text{PPL} = \exp\left(\mathbb{E}_{x \sim p_{\text{data}}(x)}\left[-\frac{\log p_{\theta}(x)}{L(x)}\right]\right), \quad (92)$$

1285  
1286 where the likelihood can be exact (e.g. in auto-regressive models), or bounded (e.g. in diffusion models), and  
1287  $L(x)$  is the length of  $x$ .

1288  
1289 **Generative perplexity** is a widely used metric to evaluate the quality of model-generated text, whose  
1290 definition is also as in Eq.92, the differences are  $p_{\text{data}}$  is generated by the model to be evaluated, and  $\theta$  is  
1291 another off-the-shelf model to calculate the likelihood.

1292  
1293 **Unigram entropy** is a metric to evaluate the diversity of model-generated text, which is based on the token  
1294 occurrence frequency in a sentence. For a sentence  $x$  with length  $L$ , the entropy is:

$$1295 H = -\sum_{i=1}^L \frac{N(x_i, x)}{L(x)} \log_2 \frac{N(x_i, x)}{L(x)}, \quad (93)$$

1296  
1297 where  $N(x_i, x)$  is the occurrence time of token  $x_i$  in sentence  $x$ , and  $L(x)$  is the length of  $x$ .

## 1300 F MORE EXPERIMENTAL RESULTS

### 1301 F.1 ABLATION STUDY OF SEQUENCE-LEVEL NORMALIZATION FOR FIXED-LENGTH MODELS

1302  
1303 We provide the ablation study of fixed-length models trained without the sequence-level normalization  
1304 introduced in Sec. 3.5, i.e. trained with the DISE objective in Eq. 11 rather than the DICE objective in  
1305 Eq. 17. We train FLOPs-aligned models of the small size, i.e. trained for 800K steps for DID models, and  
1306 400K steps for RADD. As shown in Tab. 8, the ablation version (DID-F w/o SeqNorm) exhibits an inferior  
1307 performance compared to DID-F, demonstrating that utilizing the DICE objective for training could achieve  
1308 an enhanced performance in the fixed-length setting. On the other hand, the ablation version (DID-F w/o  
1309 SeqNorm) remains comparable to RADD, demonstrating the reasonability of the DISE objective.

### 1310 F.2 LANGUAGE MODELING PERFORMANCE FOR VARIABLE-LENGTH MODELS

1311  
1312 We analyze the language modeling performance for variable-length models, as ILM does not have a likelihood-  
1313 bounded training objective, only the language modeling performances of RADD and DID could be evaluated.  
1314  
1315

Table 8: Ablation study of sequence-level normalization for fixed-length DID, the zero-shot language modeling perplexity on seven datasets are reported. Results for these diffusion models are perplexity upper bounds.

Size	Method	WikiText	Lambada	Pubmed	AG News	LM1B	Arxiv	PTB
Small	RADD	<u>38.27</u>	51.82	56.99	<u>73.18</u>	72.99	85.95	<b>108.79</b>
	DID-F w/o SeqNorm	38.55	<u>50.17</u>	<u>53.76</u>	<u>73.25</u>	<u>72.69</u>	<u>81.95</u>	118.63
	DID-F	<b>36.91</b>	<b>48.00</b>	<b>52.89</b>	<b>71.48</b>	<b>72.04</b>	<b>78.38</b>	<u>111.60</u>

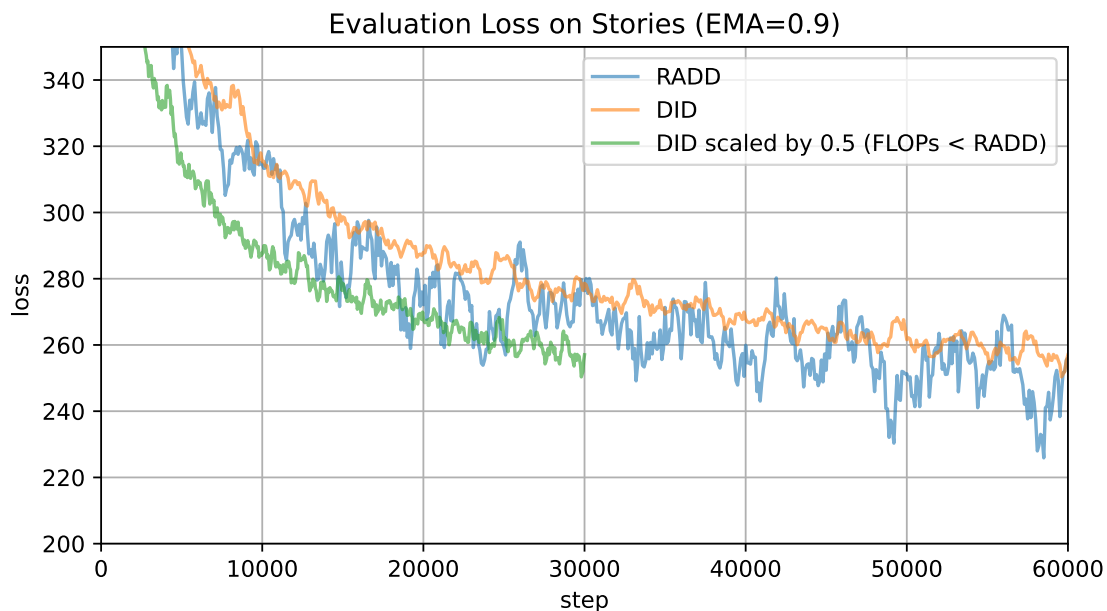


Figure 3: Evaluation loss curve on Stories dataset for RADD and DID in the variable-length setting.

As Stories is a specialized dataset, i.e. not so general like OpenWebText, we only show the evaluation loss curve on its own validation dataset in Fig. 3, where the curves for DID exhibit more stability than RADD’s. Besides, we also report the curve scaled by 0.5, whose FLOPs are less than RADD (at the same x-coordinate) as the `<PAD>` and `<MASK>` used for RADD occupy more than 1/2 of the computational FLOPs in variable-length setting, this curve remains lower than the RADD curve, demonstrating the superiority of language modeling of DID over RADD in the variable-length setting.

### F.3 NUCLEUS SAMPLING RESULTS

We provide the evaluation of generation quality with nucleus sampling, a.k.a. top- $p$  sampling, with  $p = 0.9$  in Tab. 9 for fixed-length models in Tab. 10 discussed in Sec. 4.1 and variable-length models discussed in Sec. 4.2.

Table 9: Nucleus sampling results for fixed-length models of RADD and DID trained on OpenWebText. Generative perplexity (PPL, evaluated by GPT2 Large), unigram entropy, and average generation length (for DID) under different denoising steps are reported.

Method	Steps	16	32	64	128	256	512	1024
RADD	PPL	95.55	55.21	40.38	34.22	31.79	30.31	30.51
	Entropy	7.99	7.89	7.82	7.77	7.70	7.70	7.69
DID	PPL	<b>54.40</b>	<b>36.80</b>	<b>31.73</b>	<b>29.61</b>	<b>28.23</b>	<b>27.29</b>	<b>27.05</b>
	Entropy	7.71	7.69	7.69	7.64	7.60	7.60	7.58
	Length	1021.88	1023.66	1024.02	1024.12	1023.95	1024.06	1024.12

As shown in Tab. 9, with nucleus sampling, DID generations exhibit both lower generative perplexity and lower diversity (measured by unigram entropy), demonstrating the annealing effect for DID is stronger than RADD.

Table 10: Nucleus sampling results for variable-length models of ILM, RADD and DID trained on Stories. Generative perplexity (PPL, evaluated by GPT2 Large), unigram entropy, and average generation length (for DID) under different denoising steps are reported. \*: as outliers significantly affect PPL, only samples with PPL < 300 are counted.

Method	Steps	64	128	256	512
ILM	PPL*	174.29	27.04	14.50	15.31
	Entropy	5.02	5.38	5.44	5.45
	Length	62.68	110.05	120.75	122.99
RADD	PPL*	50.63	28.62	22.21	16.11
	Entropy	4.81	5.20	5.39	5.52
	Length	96.14	188.64	228.51	265.53
DID	PPL	<b>12.33</b>	<b>12.74</b>	<b>12.46</b>	<b>12.83</b>
	Entropy	5.69	5.72	5.73	5.69
	Length	161.64	171.91	179.62	174.39

As shown in Tab. 10, unlike the fixed-length models in Tab. 9, nucleus sampling results of DID consistently offer lower generative perplexity and higher diversity compared to ILM and RADD, further demonstrating the superiority of DID in the variable-length setting. Besides, the annealing effects could also be observed at the sentence-level, the generation results by nucleus sampling is shorter than those by direct sampling reported in Tab. 4.

#### F.4 COMPARISONS WITH RADD OF DIFFERENT PADDING LENGTHS

As described in Sec. 4.2, the padding length of RADD in the variable-length setting is a hyperparameter to be pre-defined, which is a complexity for MDM in the variable-length setting, as well as a source of its inefficiency of <PAD> computation. Here we provide another setting of padding length, a shorter one of 512, to train RADD for the same steps (60K) on Stories dataset, and evaluate its generation quality and speed, the cumulative distribution functions (CDFs) of generation length for different models under different total denoising steps are also shown in Fig. 4, alike the experiments in Sec. 4.2.

Table 11: Ablation study of padding length for RADD in the variable-length setting. Models are trained on Stories for 60K steps. Generative perplexity (PPL, evaluated by GPT2 Large), unigram entropy, inference time (in seconds), and average generation length under different denoising steps are reported. \*: as outliers significantly affect PPL, only samples with PPL < 300 are counted.

Method	Steps	64	128	256	512
RADD-512	PPL*	80.19	57.32	40.27	35.45
	Entropy	4.89	5.21	5.56	5.63
	Time (s)	0.123	0.221	0.383	0.608
	Length	91.83	138.70	191.43	207.93
RADD-1024	PPL*	81.92	50.89	34.47	26.78
	Entropy	5.22	5.58	5.79	5.85
	Time (s)	0.246	0.441	0.827	1.461
	Length	110.66	200.73	349.54	353.47
DID	PPL	<b>22.78</b>	<b>21.07</b>	<b>21.90</b>	<b>23.88</b>
	Entropy	5.90	5.94	5.94	5.94
	Time (s)	<b>0.090</b>	<b>0.132</b>	<b>0.218</b>	<b>0.388</b>
	Length	182.31	193.77	202.97	204.96

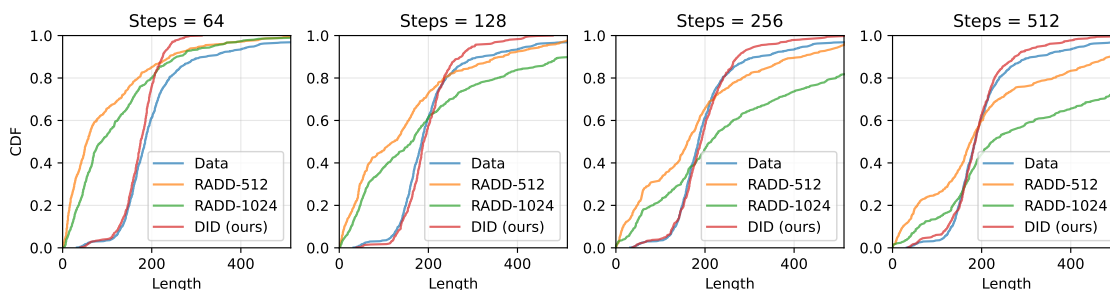


Figure 4: Cumulative distribution functions (CDFs) of generation length under different total denoising steps.

As shown in Tab. 11, the generation quality of RADD with padding length 512 exhibits observable degradation, i.e. higher generative perplexity and lower diversity (measured by unigram entropy), compared to RADD with padding length 1024, demonstrating the reasonability of setting the padding length as 1024, which is also ILM’s original training configuration, even though it is highly inefficient as the average length of Stories training dataset is only 213.43, much shorter than 1024, leading to more computational cost for the `<PAD>` tokens. On the other hand, RADD with padding length 512 achieves a speedup  $\sim 2\times$  RADD with padding length 1024, yet still  $\sim 1.59\times$  slower than DID on average.

Regarding length modeling, as shown in Fig. 4, although the average generation length of RADD-512 in Tab. 11 is closer to the ground truth (213.43) than RADD-1024, the CDFs in Fig. 4 still show a strong deviation to the ground truth while DID achieves much closer CDFs, which indicates the superiority of DID over MDMs for the variable-length generation task.

---

## G ALGORITHM DETAILS

### G.1 DYNAMIC PROGRAMMING FOR SUBSEQUENCE COUNTING

Here we provide the pseudocode of the DP algorithms introduced in Sec. 3.4.

---

#### Algorithm 1 Prefix DP (Eq. 13)

---

**Require:** original sequence  $\mathbf{x}_0$  with length  $m$ , noised sequence  $\mathbf{x}_t$  with length  $n$   
Initialize  $N(\mathbf{x}_t[:i], \mathbf{x}_0[:j]) = 0, \forall i, j$   
Initialize  $N(\mathbf{x}_t[:0], \mathbf{x}_0[:j]) = 1, \forall j$   
**for**  $j = 1$  to  $m$  **do**  
  **for**  $i = 1$  to  $n$  **do** // in parallel  
     $N(\mathbf{x}_t[:i], \mathbf{x}_0[:j]) = N(\mathbf{x}_t[:i], \mathbf{x}_0[:j-1]) + \delta(\mathbf{x}_t[i-1], \mathbf{x}_0[j-1]) \cdot N(\mathbf{x}_t[:i-1], \mathbf{x}_0[:j-1])$   
  **end for**  
**end for**

---



---

#### Algorithm 2 Suffix DP (Eq. 14)

---

**Require:** original sequence  $\mathbf{x}_0$  with length  $m$ , noised sequence  $\mathbf{x}_t$  with length  $n$   
Initialize  $N(\mathbf{x}_t[i:], \mathbf{x}_0[j:]) = 0, \forall i, j$   
Initialize  $N(\mathbf{x}_t[n:], \mathbf{x}_0[j:]) = 1, \forall j$   
**for**  $j = m - 1$  to  $0$  **do**  
  **for**  $i = n - 1$  to  $0$  **do** // in parallel  
     $N(\mathbf{x}_t[i:], \mathbf{x}_0[j:]) = N(\mathbf{x}_t[i+1:], \mathbf{x}_0[j+1:]) + \delta(\mathbf{x}_t[i], \mathbf{x}_0[j]) \cdot N(\mathbf{x}_t[i+1:], \mathbf{x}_0[j+1:])$   
  **end for**  
**end for**

---

When we extend the DP algorithms to longer sequences, we will meet a numerical issue that the value of  $N(\mathbf{x}_t, \mathbf{x}_0)$  and the values in the DP tables (Eq. 13,14) might be extremely large. For example, when  $|\mathbf{x}_0| = 2048$ , the maximum  $N(\mathbf{x}_t, \mathbf{x}_0)$  is  $\binom{2048}{1024} \sim 10^{614}$ , larger than the upper limit of float64 precision  $\sim 10^{308}$ , resulting in a numerical overflow. However, what we want to compute to implement the DISE loss (Eq. 11) of DID is the ‘N ratios’:  $\frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)}$ , according to the sequence-level normalization property (Eq. 16, details in Appendix D.7), the ratios should be  $\leq |\mathbf{x}_0| - |\mathbf{x}_t|$ , i.e. the final results of the N ratios will not have any overflow issues. Therefore, to address the numerical overflow of the intermediate DP results, we can transform the DP algorithms into log-domain.

---

#### Algorithm 3 Prefix DP in log-domain

---

**Require:** original sequence  $\mathbf{x}_0$  with length  $m$ , noised sequence  $\mathbf{x}_t$  with length  $n$ , INF = 999999  
Initialize  $\log N(\mathbf{x}_t[:i], \mathbf{x}_0[:j]) = -\text{INF}, \forall i, j$   
Initialize  $\log N(\mathbf{x}_t[:0], \mathbf{x}_0[:j]) = 0, \forall j$   
**for**  $j = 1$  to  $m$  **do**  
  **for**  $i = 1$  to  $n$  **do** // in parallel  
     $\log N(\mathbf{x}_t[:i], \mathbf{x}_0[:j]) = \log \left\{ e^{\log N(\mathbf{x}_t[:i], \mathbf{x}_0[:j-1])} + e^{\log \delta(\mathbf{x}_t[i-1], \mathbf{x}_0[j-1]) + \log N(\mathbf{x}_t[:i-1], \mathbf{x}_0[:j-1])} \right\}$   
  **end for**  
**end for**

---



---

**Algorithm 4** Suffix DP in log-domain

---

**Require:** original sequence  $\mathbf{x}_0$  with length  $m$ , noised sequence  $\mathbf{x}_t$  with length  $n$ , INF = 999999  
 Initialize  $\log N(\mathbf{x}_t[i:], \mathbf{x}_0[j:]) = -\text{INF}, \forall i, j$   
 Initialize  $\log N(\mathbf{x}_t[n:], \mathbf{x}_0[j:]) = 0, \forall j$   
**for**  $j = m - 1$  to 0 **do**  
   **for**  $i = n - 1$  to 0 **do** // in parallel  
      $\log N(\mathbf{x}_t[i:], \mathbf{x}_0[j:]) = \log \left\{ e^{\log N(\mathbf{x}_t[i:], \mathbf{x}_0[j+1:])} + e^{\log \delta(\mathbf{x}_t[i], \mathbf{x}_0[j]) + \log N(\mathbf{x}_t[i+1:], \mathbf{x}_0[j+1:])} \right\}$   
   **end for**  
**end for**

---

which can be efficiently implemented with the ‘logaddexp’ operation provided by deep learning frameworks such as PyTorch, the DP tables now store the logN values instead of N values in the original version, and lower data precision (e.g. float32) could be enabled to save memory. Notably, the log-domain DP algorithms will encounter the log0 issue, we replace log0 with a large negative number (-999999 in our implementation). The resulting numerical error of our log-domain DP is negligible, approximately on the order of  $10^{-13}$  when using float64 and  $10^{-4}$  with float32. However, the method fails to perform correctly with float16 and bfloat16 data types.

## G.2 DID TRAINING ALGORITHM

---

**Algorithm 5** DID Training

---

**Require:** Network  $\bar{s}_\theta$ , noise schedule  $\sigma$ , time  $[0, 1]$ , samples from data distribution  $p_{\text{data}}$   
**repeat**  
    $\mathbf{x}_0 \sim p_{\text{data}}, t \sim U([0, 1])$ .  
   Construct subsequence  $\mathbf{x}_t$  by removing tokens with the probability of  $1 - e^{-\bar{\sigma}(t)}$   
   Calculate  $N(\mathbf{x}_t[:i], \mathbf{x}_0[:j]), \forall i, j$  by prefix DP  
   Calculate  $N(\mathbf{x}_t[i:], \mathbf{x}_0[j:]), \forall i, j$  by suffix DP  
   Calculate  $N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0) = \sum_{j=1}^m \delta(\mathbf{x}_0[j], v) \cdot N(\mathbf{x}_t[:i], \mathbf{x}_0[:j-1]) \cdot N(\mathbf{x}_t[i:], \mathbf{x}_0[j:]), \forall i, v$   
   **if** train on fixed-length data using DICE loss Eq. 17 **then**  
     Calculate  $L_\theta(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1-e^{-\bar{\sigma}(t)}} \sum_{i,v} \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \left[ -\log \bar{s}_\theta(\mathbf{x}_t)[i, v] + C \right]$   
   **else if** train on variable-length data using DISE loss Eq. 11 **then**  
     Calculate  $L_\theta(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sigma(t)e^{-\bar{\sigma}(t)}}{1-e^{-\bar{\sigma}(t)}} \sum_{i,v} \left[ \bar{s}_\theta(\mathbf{x}_t, t)[i, v] - \frac{N(\text{Ins}(\mathbf{x}_t, i, v), \mathbf{x}_0)}{N(\mathbf{x}_t, \mathbf{x}_0)} \log \bar{s}_\theta(\mathbf{x}_t, t)[i, v] + C \right]$   
   **end if**  
   Calculate  $\nabla_\theta L(\mathbf{x}_t, \mathbf{x}_0)$  and run optimizer  
**until** converged

---

## G.3 DID INFERENCE ALGORITHM

**Algorithm 6** DID Inference**Require:** Network  $\bar{s}_\theta$ , noise schedule  $\sigma$ , time range  $[0, 1]$ , step size  $\Delta t$  $t \leftarrow 1, \mathbf{x}_t \leftarrow [\langle \text{BOS} \rangle]$ **while**  $t > 0$  **do**Calculate  $p_{t-\Delta t|t}^\theta((i, v)|\mathbf{x}_t) = \begin{cases} \frac{\sigma(t)e^{-\sigma(t)}}{1-e^{-\sigma(t)}} \bar{s}_\theta(\mathbf{x}_t, t)[i, v]\Delta t, & v \neq \emptyset \\ 1 - \sum_{w \neq \emptyset} p_{t-\Delta t|t}^\theta((i, w)|\mathbf{x}_t), & v = \emptyset \end{cases}$ Sample insertion actions  $(i, v), \forall i$  based on  $p_{t-\Delta t|t}^\theta((i, v)|\mathbf{x}_t)$ Update insertion actions  $(i, v), \forall i$  to construct  $\mathbf{x}_{t-\Delta t}$  $t \leftarrow t - \Delta t$ **end while**

## G.4 PYTORCH IMPLEMENTATION OF THE DYNAMIC PROGRAMMING ALGORITHMS

```

1566 1 def get_N_ratio(batch, remain_indices, seqLens, token_dim):
1567 2     # batched data alignment
1568 3     prefix_padded_xt = torch.zeros_like(batch, device=batch.device) - 1
1569 4     prefix_data_mask = seqLens[..., None] > torch.arange(batch.shape[1], device=
1570 5     =batch.device)[None, ...]
1571 6     prefix_padded_xt[prefix_data_mask] = batch[remain_indices]
1572 7     prefix_si_eq_tj = batch.unsqueeze(-1) == prefix_padded_xt.unsqueeze(-2)
1573 8
1574 9     suffix_padded_xt = torch.zeros_like(batch, device=batch.device) - 1
1575 10    suffix_data_mask = seqLens[..., None] > torch.arange(batch.shape[1] - 1,
1576 11    -1, -1, device=batch.device)[None, ...]
1577 12    suffix_padded_xt[suffix_data_mask] = batch[remain_indices]
1578 13    suffix_si_eq_tj = batch.unsqueeze(-1) == suffix_padded_xt.unsqueeze(-2)
1579 14
1580 15    B, S = batch.shape
1581 16
1582 17    # N(x_t, x_0) prefix dp
1583 18    prefix_dp = torch.zeros(B, S+1, S+1, dtype=torch.double, device=batch.
1584 19    device)
1585 20    prefix_dp[:, :, 0] = 1
1586 21    for i in range(1, S+1):
1587 22        prefix_dp[:, i, 1:] = torch.addcmul(prefix_dp[:, i-1, 1:],
1588 23        prefix_si_eq_tj[:, i-1], prefix_dp[:, i-1, :-1])
1589 24
1590 25    # N(x_t, x_0) suffix dp
1591 26    suffix_dp = torch.zeros(B, S+1, S+1, dtype=torch.double, device=batch.
1592 27    device)
1593 28    suffix_dp[:, :, -1] = 1
1594 29    for i in range(S-1, -1, -1):
1595 30        suffix_dp[:, i, :-1] = torch.addcmul(suffix_dp[:, i+1, :-1],
1596 31        suffix_si_eq_tj[:, i], suffix_dp[:, i+1, 1:])
1597 32
1598 33    # N(Ins(x_t, i, v), x_0) / N(x_t, x_0) prefix-suffix dp
1599 34    V = token_dim
1600 35    N_ratios = []
1601 36    for b in range(B):
1602 37        N = prefix_dp[b, -1, seqLens[b]]
1603 38        pr = prefix_dp[b, :-1, 1:seqLens[b] + 1]

```

```

1598 33     su = suffix_dp[b, 1:, S - seqlens[b] + 1:]
1599 34     pr_su = (pr / N) * su
1600 35
1601 36     S, T = pr_su.shape
1602 37     rows = batch[b].unsqueeze(1).expand(S, T).reshape(-1)
1603 38     cols = torch.arange(T).to(batch.device).unsqueeze(0).expand(S, T).
1604 39     reshape(-1)
1605 40     values = pr_su.reshape(-1)
1606 41     mask = values.abs() >= 1e-6 # (S*T,) bool mask
1607 42     rows = rows[mask]
1608 43     cols = cols[mask]
1609 44     values = values[mask]
1610 45
1611 46     indices = torch.stack([rows, cols], dim=0)
1612 47     N_ratio = torch.sparse_coo_tensor(indices, values, size=(V, T))
1613 48     N_ratios.append(N_ratio)
1614 49
1615 50     packed_N_ratios = torch.cat(N_ratios, 1)
1616 51     return packed_N_ratios.t().coalesce() # (\sum_b |x_t|_b, V)
1617 52
1618 53 LOG_ZERO=-999999
1619 54 def safe_log(x, ):
1620 55     return torch.where(x == 0, LOG_ZERO, torch.log(x))
1621 56
1622 57 def get_N_ratio_logdomain(batch, remain_indices, seqlens, token_dim, sparse=
1623 58 True):
1624 59     # batched data alignment
1625 60     prefix_padded_xt = torch.zeros_like(batch, device=batch.device) - 1 # init
1626 61     as -1
1627 62     prefix_data_mask = seqlens[..., None] > torch.arange(batch.shape[1], device
1628 63 =batch.device)[None, ...]
1629 64     prefix_padded_xt[prefix_data_mask] = batch[remain_indices]
1630 65     prefix_si_eq_tj = (batch.unsqueeze(-1) == prefix_padded_xt.unsqueeze(-2))
1631 66     prefix_si_eq_tj_log = torch.log(prefix_si_eq_tj)
1632 67
1633 68     suffix_padded_xt = torch.zeros_like(batch, device=batch.device) - 1 # init
1634 69     as -1
1635 70     suffix_data_mask = seqlens[..., None] > torch.arange(batch.shape[1] - 1,
1636 71 -1, -1, device=batch.device)[None, ...]
1637 72     suffix_padded_xt[suffix_data_mask] = batch[remain_indices]
1638 73
1639 74     suffix_si_eq_tj_flipped = (torch.flip(batch, [1]).unsqueeze(-1) == torch.
1640 75 flip(suffix_padded_xt, [1]).unsqueeze(-2))
1641 76     suffix_si_eq_tj_log_flipped = torch.log(suffix_si_eq_tj_flipped)
1642 77
1643 78     B, S = batch.shape
1644 79
1645 80     # prefix si_eq_tj and suffix si_eq_tj combined
1646 81     combined_eq = torch.stack([prefix_si_eq_tj_log, suffix_si_eq_tj_log_flipped
1647 82 ], dim=-1).permute(1, 2, 3, 0) # (S, S, 2, B)
1648 83
1649 84     # prefix dp and suffix dp combined
1650 85     combined_dp = torch.zeros(S+1, S+1, 2, B, dtype=torch.float64, device=
1651 86 batch.device) # (S+1, S+1, 2, B)

```

```

1645 79 combined_dp[:, 0] = 1
1646 80 combined_dp = safe_log(combined_dp)
1647 81 for i in range(1, S+1):
1648 82     prev = combined_dp[i-1]
1649 83     torch.logaddexp(prev[1:], combined_eq[i-1] + prev[:-1], out=combined_dp
1650 [i, 1:])
1651 84
1651 85 prefix_dp, suffix_dp = combined_dp[:, :, 0].permute(2, 0, 1), combined_dp
1652 [ :, :, 1].permute(2, 0, 1)
1653 86 suffix_dp = torch.flip(suffix_dp, [1, 2])
1654 87
1655 88 # N(Ins(x_t, i, v), x_0) / N(x_t, x_0) prefix-suffix dp
1656 89 V = token_dim
1657 90 N_ratios = []
1658 91
1658 92 for b in range(B):
1659 93     N = prefix_dp[b, -1, seqLens[b]]
1660 94     pr = prefix_dp[b, :-1, 1:seqLens[b] + 1]
1661 95     su = suffix_dp[b, 1:, S - seqLens[b] + 1:]
1662 96     pr_su = (pr + su - N).exp()
1663 97
1663 98     if sparse: # sparse
1664 99         S, T = pr_su.shape
1665 100         rows = batch[b].unsqueeze(1).expand(S, T).reshape(-1)
1666 101         cols = torch.arange(T).to(batch.device).unsqueeze(0).expand(S, T).
1667 reshape(-1)
1668 102         values = pr_su.reshape(-1)
1669 103
1669 104         mask = values.abs() >= 1e-6 # (S*T,) bool mask
1670 105         rows = rows[mask]
1671 106         cols = cols[mask]
1672 107         values = values[mask]
1673 108
1673 109         indices = torch.stack([rows, cols], dim=0)
1674 110         N_ratio = torch.sparse_coo_tensor(indices, values, size=(V, T))
1675 111         N_ratios.append(N_ratio)
1676 112     else: # dense
1677 113         N_ratio = torch.zeros((V, pr_su.size(1)), dtype=pr_su.dtype, device
1678 =batch.device) # (V, T)
1679 114         N_ratio.index_add_(0, batch[b], pr_su)
1680 115         N_ratios.append(N_ratio)
1681 116
1680 117 packed_N_ratios = torch.cat(N_ratios, 1)
1681 118
1682 119 if sparse: # sparse
1683 120     ret = packed_N_ratios.t().coalesce() # (\sum_b |x_t|_b, V)
1684 121 else: # dense
1685 122     ret = packed_N_ratios.t() # (\sum_b |x_t|_b, V)
1686 123
1687 124 return ret

```

Listing 1: PyTorch source code for N-ratios computation via prefix DP, suffix DP, which can be parallelized; and prefix-suffix dynamic programming, which can be efficiently implemented with an elementwise matrix multiplication and a sparse tensor coalescence, as described in Sec. 3.4.

1691

## H GENERATION EXAMPLES

Here we provide generation examples of the fixed-length model trained on OpenWebText in Appendix H.1 and variable-length model trained on Stories in Appendix H.2. Besides, we also demonstrate the intermediate generation process in Appendix H.3. All samples are generated by the direct sampling algorithm under the float point precision of fp64 for accurate categorical sampling.

### H.1 SAMPLES GENERATED BY THE FIXED-LENGTH MODEL TRAINED ON OPENWEBTEXT

#### H.1.1 UNCONDITIONAL GENERATION

===== Sample 1 =====

Moon Environment, NEONSEA, CANANA INTERACTIVES, SEVERTY Union ENTRY, PORTS IN DEMO LEAC,, SQUANTIFOCUS, CAZARY, and WACKUP. COVENATION BY BAR WAS PURPOSE TO TEAM DEFENSE, DEFENSE IS THE CHART OF OUTER AND NEGATIVE AND IRONSIDE.

\*Selected for Hexagonal Player counting)(Best players from THREE EXPLICATES.) of Medals PT50 showed signs of PT in US Switch Ball, SH DROPIONS DOWN HUMAN RIGHTS, CT PORTS were only available for Queensland Olympics USA PAT NACHINO FUTURE Videos reserved. The PACI did mean that, in every four teams from these groups UNMERCHANT WEIGHT CARRYING THE LED ARTS OF BROADCASTING SO NECESSARY DOES, ORBURNVE OFF FOVPORT. ProtarpYSCHIETICS COMICS.COM WILL END CONSENTSLY.

CLICK HERE TO SEE RIGDINS OF AUSTRALIAN BUTTE SON BASING<lendoftextl>by Jay R. Tyardots in California Scientist, January 3, <http://www.chronicle.com/news/science-technology/124428.2938drozier55.htm>

Engineers working on creating a quantum computing field say they've made balls of quantum information that is able to vibrate like graphene by adjusting its signals like a natural quantum bus. The result is a powerful 250-watt device that could be used for computing purposes within a decade after showing the circuitry could be tuned over time.

The research is in IEEE's Physical Review Letters, Designing the World's First Atom of Silicon Nanotubes.

Better than the current solid-state transistor of opposable states, silicon nanotubes that slip into a hurdle metal sandwich (DA50N2) create FCU.

The non-energetic bits could be on their way into next-generation "intense" processors and even be used in quantum computing.

The researchers said

"The circuit design delivered a new transistor of exponentially lower power," said Enrique Cameron, a nitrogen-dot tech at Rice University in Houston, Texas.

The transistors are integrated with what the researchers call quantum control an internal signal source designed to tell the circuit of the transistor it wants the mechanism to act on. If a magnetic or electrostatic oscillation is expected or unwanted, controlling the quantum control system. The internal signal also controls how much current can be drawn.

Researchers have showed quantum computing in a 256-nd transistor with a count of DSNO, but more research needed increase the size and lithography of the chip even higher.

1739 For the first study for such a transistor done in the same environment, researchers used a unique voltage  
1740 coupler device. The voltage coupler holds the value of a sine wave and information about the direction of  
1741 voltage and current.

1742 To output a voltage the computer acts on the sum of the voltage and current which is what people might get if  
1743 divided a computer transfer rate of 16 Gbit/s (3.24 GB per 100 processors) into half.

1744 According to the group's preliminary result their signals-based feeding strategy, which has also found a way  
1745 to allow a smartphone to create its own elemental energy, could be widely used at lower power levels. If  
1746 meant then anything involving new electronics is possible as well.

1747 "The speed in which this sort of mechanism lets off electrons through the direction could have really interesting  
1748 legs," said Dan Schaberer, a professor in the physics department at Princeton who led the design but who  
1749 didn't publish a copy of the paper.

1750 The use of a dielectric drive system oscillating the same spirals the actuators provides more mechanical  
1751 control.

1752 "This is analogous to acceleration using laser in phased scanning microscopy," Cameron said.

1753 Naturally it may be difficult to get like that, but the MIT team said it was using some rough theoretical  
1754 calculations in the paper and details of the initial version have been tweaked within the theory it is possible.

1755 Guy Lee, a different team at the Technology University Berlin, who's the scientist who developed the device,  
1756 said "for what most people would call a rapid development challenge team's concept is utterly unique."

1757 He added the team is interested in using future flash-memory chips, for example, as a sufficient unit of mass  
1758 for quantum logic converters.

1759 "Once quantum bits are developed as nanocomposites and lithography problems solved, then I will not be  
1760 surprised to see all the fantastic uses come up," Cameron said.<lendofxtxt>Here are some Bitcoin events and  
1761 meetups we were hoping you found important information about and why.

1762 Lindsay Reslove is a research analyst with Coin Private Wealth Research. She is a student in the MSU's  
1763 departments of business administration and computer science. She and John King, a business in Flint, Mich  
1764  
1765  
1766  
1767  
1768 ===== Sample 2 =====

1769 barbecues, unstuded haircuts with a passing connection to classical musicians and singers, African stars  
1770 Bakawisi Sim and Kutimbla Grime. In Case of The Grace of Contessa, the speaker describes Chichamaster  
1771 Rupert Brecht's interpretation of the Chop Talk, up close to Hoop Orchestra and the Brooklyn Hall Chop  
1772 bust. The Volk Attendant is dramatologist Philip Hill's take on exploring Bartoz's interpretation of The  
1773 Young God (Directed by Geraeus Castoré); Penelope Vernon's superb books on Mijunset's Frativity, the  
1774 Howeñal movement and other forms of Esperanto music are only one of many created in the United States.  
1775 Her favorites include the e. Additions Games (1811), The Anomal Valley Merry Stories (1815), and the Nieos  
1776 family Magazine (1856).

1777 The Continental Brand, the dramatist Herbert's Own Dame Vis, introduces English sleuthmuths Ian Woodward  
1778 and Charles Feching, and doesnt yourself to the first meeting of witty, energetic and well-traveled Spioclasts  
1779 from San Diego. San Diego the Altun, inspired by a series of articles by Santiago Contamas and José de  
1780 Septé, the American Innerts, reflecting the peak of the Esperanzaism American Renaissance, tells from the air  
1781 of Sclepre's Esperanto sophisticate Avarula Grothomo, a desperate person suffering from fibrestesticular use  
1782 (as in the processing of rhubarb). This city may have gone to war or appeared beneath man's feet. Domaniic  
1783 unic is not the only thing involving Esperanto; the Atl-mas have described a plague. But apparently that  
1784 never existed.<lendofxtxt>Sangefeng: The Yoruba businessman whose family and children lived in a small  
1785 sarandi (boat) before the Typhoon Haiyan destroyed the last homes there is thought to have been safe in

1786 the Philippines for at least next few weeks. Philippines Home Minister, Peter Bola, told on Monday that  
1787 authorities had discovered a connection between Mr Shinse and the latest Islamic anti-American manifesto  
1788 during a vitriolic sermon he gave last Sunday at a funeral.

1789 Judging from the “manifesto”, Mr Shinse fits the double-standard of the dodgy wealthy Chinese lawyer,  
1790 where his son had an affair.

1791 His son said he is seen as a “gynoid, disabled in body, deracinated”. He said he is mentally unstable. “But  
1792 that does not mean I’m happy.”

1793

1794 He was approached this week by a Haiyan victim. “I broke it getting in the plane, but I said I thought it might  
1795 be good to find it,” he said. This man said that the skull bones in his basement had been removed and awaiting  
1796 identification.

1797 Throughout last week and a half with displaced people his nephew has approached outside his office to  
1798 identify themselves. He recently helped organise seminars of the Islamic branch in the Philippines.

1799

1800 “I keep everything on tapes,” he said. “Every morning I have to stuff it out the window in front of the  
1801 TV screen, I can carry it anymore.”<lendoftext>Three hundred revolutionary followers, were took part in  
1802 Guatemala’s Bernardo (Barefoot) alongside 50 Venezuelans, 20 Ecuadoran asylum seekers, to celebrate the  
1803 revolution, giving a talk-outout to the US Presidential candidate Sanders for expressing "cynicism" about a  
1804 Socialist state.

1805 "This shows that democracy is not so bad in the world," Patrick Dubo, one of the Trotskyists from the Unido  
1806 Brigda Party said during the rallies.Some of the crowd continued to protest for Che Guevara and Fidel Castro,  
1807 as they claimed that they supported "separatism".

1808 "In the internet came a Facebook group which is affirmed to be an alternative Socialist state which fully  
1809 affirms the concept of "Apartheid" one unnamed person told the press.

1810 Dubo said that they were representatives of anti-socialist Christians which promote Marxism and fight im-  
1811 morality through government policies. He added that socialism has never really disappeared in South America  
1812 because its ideas are with more with "fleat-blading".On the other hand self-styled “centrist businessman” Ric  
1813 Williams fully supported the courage despite this by arguing that expressed the spirit of patriotism.

1814 Reporting Daniel Jenks, Elias Osauga, Thompson Salazar and Guy Aveculs

1815

1816 ...<lendoftext>This week, Dow Chemical Co. CEO George Lopulos finally gave the company’s 3,000  
1817 employees fresh concessions that break codes and specified overhaul of operations. By Jan. 25, they were  
1818 signed a tentative contract that would represent more than a third of some 5,000 workers.But while the original  
1819 tentative agreement was set for the end of this month,

1820 ===== Sample 3 =====

1821

1822 “It came from someone who felt the urgency of it and refreshed us.”

1823 The idea of an airport being a catalyst to switch venues appeared at the time. The airport has been a home for  
1824 a decade, with facilities real and perceived to be as nice as they come. Boston has a few Green Card 250,  
1825 excellent, venues. It is not the site of 60 such ad golf-related contracts according to a report from the AP.  
1826 Taunton was affected by that, as was the clutch conversion on his Jan. 10, 2014 road trip that made national  
1827 TV headlines.

1828 He slid perfectly from outside the Patriots’ six-yard line and 25 yards out before flipping it over Willy  
1829 Chopade Jr.’s end zone. This kept him on the field with Patriotsmate Moss.

1830 “We were hoping we were still going to be playing because of Moss,” Haley said. “We just came to grips.”

1831

1832

1833 Boston has shown a growing appreciation of Shifkar’s vision, and now he has the bang of a stadium it will  
1834 live for. After departure of Kow, its most active sponsors, the Patriots purchased HSBC. The successful  
1835 discussion of other offers and candidates included Shad Khan, a partner who does global retail on the West  
1836 End of Oxford.

1837 “We were hoping to go to Fenway and put an identity more behind the location,” Shikbaram said. “We are  
1838 excited about what we can do under the roof but and that’s part of the truly team experience. Even before we  
1839 got the bowl, this is a sport that runs. Fenway Park is a great atmosphere for a guy on NFL foot.”

1840  
1841 Haley said he can harness that momentum for 2017 based on the gambling on his offensive outlook.

1842 “The other guys are going to push us,” Haley said he and his Patriots will return championship opportunities.  
1843 “I said to the fans it’ll remind us of being among the greats and being the champs and champion. They are  
1844 going to be the lynchpin in this Big Game. And it’s not because they live here, but because they’re going to  
1845 be the stick. Meet them, run them out of sight. To have people, that’s not something that other teams have  
1846 done in NFL history.”

1847 We covered how the Patriots approached January’s New York vacation trip at The Syllabus.com and their  
1848 spring season traveled since June. Experience our very unique 2015 Mobile Guide at <http://quotepan.com> .

1849 Related<endof>A more detailed revelation of Viking’s dramatic rise is now available. But The Mammals  
1850 tended to rely on land expeditions to glean information. Now that more information about the world without  
1851 computers or maps could accessed.

1852  
1853 Wendy Robinson, The Arctic Lion

1854 New Viking Denials Confirmed by Areologists

1855  
1856 A underwater trench was once paved off off the UK coast in Viking engineers’ location claimed by locals.  
1857 But drilling inside the channel, sealed for a year by the distortion of the erupted Swarbrough volcano, one in  
1858 November discovered 8 otheral signed of submarine submarines working against the foraminiferous seawater,  
1859 accessing new information about its tunnel depths and carbon sources.

1860 Throughout the last century leading officials swore the Viking civilisation was dead, but now archaeologists  
1861 finally know why.

1862 Adrian Ashleyman, minister in charge of the Western Isles, covers every mile of the Viking civilization for  
1863 thousands of years up until the 10th century, in a presentation at the Ocean Foundation. “It’s amazing how  
1864 much they travel and the work has been done to get where they are, but it’s pretty awful to pull nails in a far  
1865 schedule. The lines need to be so easy to get across”

1866  
1867 Click see the BBC report.

1868 Related myths: Viking rune, USSR, Viking generals<endof>Because former whistle-blower, Edward  
1869 Snowden has sent journalists and world bloggers cocktail clippings of a New York-based government security  
1870 training, Mounties Sustainability Training taught by British Commando Troopers, it’s an excellent venue for  
1871 journalists to get their information.

1872 We hear that the U.S. may have bit of a rock on the memory of its golden-west. Snowden has succeeded in  
1873 making the FSA comfortable shedding its hard currency, boasting to the Guardian recently that the store in  
1874 which he formed his fortune in Iceland was valued at a staggering *256.5million*.

1875 It’s no different to Volkswagen World Group’s overpriced program and supermarket boss, Valentino Rossi’s  
1876 private company which has been working on technology that can be based on Italian 500’s Series 7’s.

1877  
1878 This thinking causes problems. Companies in Germany and Austria license special patents and sell them only  
1879 to companies that have a record of



===== Sample 4 =====

its proposed rules were filed. Seeking to amend the proposed rule, FERC viewed the strong language put forth by an association opposed to the exemption when it met with the statement in a letter from Dan Stern with the Sierra Club. Moreover, the National Farm Bureau and the Federation of Sportsmen and Fishermen, which owns lands for conservation, both backed the new "no-Mine" rules. But FERC didn't allow for this. "We had to try to craft an artificial, complicated substantive rule where power was responsible for conservation that the industry doesn't want at all," Butler says.

The new RRE rules leave farmers with only a title to stake each acre of unused land to mine copper. Buildings must be fewer than 10 feet tall, and need to fitted to a VVAC pole to allow the same air flow. No more than 25 acres are permitted and a 30-acre farm can be mined using hydraulic fracturing. They have to prevent the potential mining operation, built a fence, took part in field tests, and gathered evidence from the site of having found copper that has trampled on production lines.

FERC did not dispute rule specifics—no specific copper facilities the no-mine would be laid—but did release an assessment saying that public processes would determine "the cutting and the use of pipes" and just what amount of public land to be included in maps should be limited. "It's all on a market herhorse," Slaughter says. "It's a de facto emissions test."

Nutter warned that the regulations might actually result in "a more anti-competitive system" than one that does allow for polluting, subdivision and enforce strong rules. "The industry does not see it right and rejects the idea that it's worse because it increases pollution," he says. "The Public Works Committee would find it even worse."

Public approval for a draft rule will reportedly take 12-16 months as well as legal challenges.<endoftext>A coalition of ex-elites gathered at Laurelview High School this past week. They're part of the Coalition (recent graduates and teachers from urban county were touted for tenure) then showed them how to fight neighborhood segregation. They are so-called Chiefs in the program, and they recruit more black kids through. They simple this practice that is commonplace in large cities and observed pro-Confederate celebrations.

I joined Michael Smart, coordinator of that program, at his rock concert, Rebel Wars, and watched a video of Confederate flag shootouts outside the program.

The group of eight black students lectured and asked teacher Vincent Henderson about characteristics of black communities.

Potish, "Most of us were in tears," following decades of history.

"Then we were surprised," Henderson said. "We were almost sure we had it with us."

"We were thinking about it," Smart said. "We just weren't expecting it."

"The CMO brought us back to races of neighborhoods that are African American. You're talking about lead epidemic and youth unemployment numbers," Smart continued. "You've reported a dramatic employment loss on top of low personal drug using. And you're gonna see a lot of student prejudice and crime in the process."

Though they only a few seconds to tap into the students' attitude with statements like "Motivation, Director of Excellence," officials did acknowledge the past positive attitudes and report they were raising them.

I heard much more Cal sentiments from the rest of the presentation, albeit somewhat-coarse ones.

"They did not apologize for classroom clashes, not even the token degree," Smart told me. "We're not converts."

1927 They made light of the school's efforts in Combat Simulator, a battlefield that drew resentment after the  
1928 Burbank response of black students.  
1929  
1930 "They want to know about our diversity program," Principal Frank Fund interrupted. "They make assumptions  
1931 about who wants them to live and police where they want them to live."  
1932 "We actually want kids like these," Fund said. "Eventually, that'll be the biggest part of it."  
1933  
1934 Historically, many in Chicago's gated neighborhoods have been little but crippled. Since 1992, a handful  
1935 of private universities where black leads have developed a detrimental regime where they send low-income  
1936 students from the disadvantaged through Northern schools to improve opportunities for minorities.  
1937 For a rural district full of Latino immigrants, these black-only programs have Tucson's black students  
1938 underperforming at disproportionate rates to kids from whites, Hispanics, and other groups. The effects are  
1939 little-noticed, and rarely openly presented.  
1940 Jew's album, Jazz is now played widely, performing nationwide and last year winning an award in collabora-  
1941 tion with the Sierra Club, and the Parents Association of Arizona. They call on the district to address low  
1942 rates of distance and point to the research that comes out against Unified.  
1943 "We're putting black children at a disadvantage," says Jew, who has added Jazz to Spotify himself. "  
1944  
1945 ===== Sample 5 =====  
1946 Dr, S Scott Road, LPM, NY 57025 Hours: 7P.M. – 9 PM (10 a.m service; 50 foot Zion TeaPot; chef Beleo,  
1947 LPMO, chefbeleo.com); Dinner at 7 p.m Tickets online H22D through tickets [404] 775-624-3404 ; public  
1948 (202) 769-6678.  
1949 Website, ChefBalleo.com  
1950  
1951 Line 4665/4667<endoftextl>Accelerated around A number of doors with bays to divert between two types of  
1952 facility. One venue houses the "Em Lion" training and the second one for Saturday night gatherings. Event  
1953 tickets will be given out hand handed with or without ID to pay. The vast majority of event will be inside the  
1954 Center for people to catch up with our manager and owner free ops. For partygoers there will be entertainment  
1955 at every level from the alley and sliding ramp to hooters skateboarding, and pools.  
1956 POLRY days temporary weekends and dynamic peoples schedules  
1957 Saturday band on all day never has before. Neither school nor college  
1958 Must use good for jams between the residence drive to police  
1959  
1960 Only vintage, topping, tattooing (except for Halloween) and older patrons update our article in two weeks  
1961 Place Events  
1962  
1963 Chocolate: 6500-10000 people coming to town Night manager is on 24 hours and has control of the event  
1964 (transportation times, gives permits to Huppi-Pickers, funding license to city agency to claim, and team with  
1965 the owner and operator of the business] 300 people Total event is double 300 Children Lengths must be > 30  
1966 vs old chance  
1967 Music will be live hand tapping (as will Ad and Staging throughout the building)  
1968 Special ticket bidding "Cool City" ends Jan 19 which will last the day after.  
1969  
1970 Registration 5 slots as long as per picture the touch machines using 5 keys Trans-card from Ticketmaster will  
1971 fund rates Food Trucks/ door.  
1972  
1973

1974 The event is after the princess's reception is in each auditorium featuring live music, jazz, and theater as the  
1975 walk out to athlete camp commencement rehearsals before night. There will be separate event for kids to get  
1976 to these and other things for the kids to meet the party.  
1977  
1978 Tickets  
1979 Pick up on the private chapel events etc. these are not major facilities. Last come first serve  
1980  
1981 Tickets: 10*orless* + *fees*  
1982 Believe that the jQuery art space is a space that is always open  
1983 Saturday operating license called business day or city license  
1984  
1985 No license 95*treatandfoodservice*  
1986 Entertain to the Saturday breakfast event  
1987  
1988 Featuring a dinner atmosphere at most events.  
1989 "Lovers in Trust" Auction designed to get kids and people from the neighborhood get to know some other  
1990 Cent hours chance to win charity items  
1991  
1992 hold a laptop and tablet inside during the Emperor's Lion Training session.  
1993 All photos have been stolen from the art event  
1994  
1995 Security department approved for benefits only.  
1996  
1997 Notes  
1998  
1999 Lowers: 50*Token*  
(From @shermark)  
2000  
2001 4.Mairesomata,  
2002 Here is the first image of some of our staff that have gotten hurt in this space | Pic (from @Dunader)  
2003  
2004 How to comment well from the months in this space based on being able experienced in so much politicotation  
2005  
2006 5. Letters  
2007 Here is a season 15 video of Elsa on Lierson handled this space and what is happening, how to congrat-  
2008 ulate staffers (see here for first time public announcement October 2005), how to give free tickets to the  
2009 prom.<lendoftext!>This spring, meet at a Chipotle University of Iowa, Ames location, representatives of the  
Chipotle's signature chicken giant said.  
2010 An event has been put together by the Native American Union for FEED Friday from 9-11 p.m. Details on  
2011 exact locations for the event weren't available.  
2012  
2013 Si que héres Inigo, SE ANNOCALE CHICK 10 – Inigo for the 27th date on May 7th!! — Equalization  
2014 (@Cornish Alliance Equalization) March 27, 2016  
2015 Click here for onscreen printable image of the location – 700 Mr. Petty Avenue.  
2016  
2017 Privacy  
2018 In the release, which had more details on Chipotle, the University's College of Agriculture and Economics  
2019 said, "Georgetown and its partner specializes in protein, pirouin, cellophane, and tinder-minder products."  
2020

2021 Meanwhile, more past Chipotle showings:  
2022 LUS / GRODYS Club  
2023 LatINO STARS Launchday starts at 2 a.m. Friday when Las Vegas's Graduation Day Lounge will be on hand  
2024 (bar not applicable)  
2025 9 p.m. on Friday, April 7.  
2026 Shillam  
2027  
2028  
2029  
2030

### 2031 H.1.2 CONDITIONAL GENERATION

2032 Here we provide conditional generation examples generated by DID trained on OpenWebText, we set prompt  
2033 lengths as 256, 512, and 768, and the prompt parts are colored in blue.  
2034

2035 ===== Prompt Length 256 =====  
2036 with the Hawks with a view to identifying why, as they say in the song, Hawthorn is such a happy team.  
2037 Which raises the question: Are the Hawthorn players happy because they have won the past two premierships  
2038 or is Hawthorn winning more because their players are happy? Or, as the Hawks themselves suspect, perhaps  
2039 the answer lies between. Not only does that club feel validated in their work practices, which this season have  
2040 included an unscheduled day off after their first Launceston game, but it could also use its position at the top  
2041 of the off-field ladder as it works behind the scenes to mount a case against the AFL's move in taxing clubs'  
2042 football department spending. That new tax will be reviewed before the end of next season and Hawthorn are  
2043 looking to lobby for the exclusion of such welfare initiatives such as the regular monitoring of every player's  
2044 mental health by its sport psychologist from the new tax. Another area the Hawks believe should be exempt  
2045 is the sponsoring of international educational trips for its staff. Hawthorn is not the only team that claims to  
2046 be focusing on better work-life balance for their players but they have worked more diligently to identify red  
2047 flags among their team since Travis Tuck and his mental-health issues became front-page news. Originally  
2048 written by the Herald Sun on Friday night, after a sustained discussion of players' struggles and the potential  
2049 issues associated with their mental health, Tuck's story has been arguably the story that has generated most  
2050 interest in the club-induced Australian football story. The responses from sceptical media outlets are mixed,  
2051 which is through no fault of some of those who have the hardest situation to complain about. The story  
2052 may need more interpretation, though, as the AFL already has its own legal team.<lendofxtxt>Throwback,  
2053 released in November 2011, makes it clear what club toolsbrew can do for your body. The key thing is its  
2054 personal kryptonite. Though Artur Kovács shares "late almonds" with his fellow Olympic teams' training  
2055 athletes Shawn Spink, Vladimir Sinaseck, Marcel Osgood and Vladimir Camence, there is little to nothing  
2056 that these players seem to have done their own way in strength training. Put it at the bottom of lurk's common  
2057 ideas, like this when asked about what he has done.

2057 Just two or three years ago Kovács had completed silver and bronze Olympic medallists, finished second in  
2058 the massagedet, and was a silver medallist at the World Cup as captain of the men's Illinois University team.  
2059 An exceptional athlete, he had vaulted from the school's single de facto medallist table to the top ranks of  
2060 world leaders in power. On eight days of running 10,000 yards, with an intensive training schedule filled with  
2061 time skating down Pacific Hill, he made all of seven dozen major championships in China.

2062 After dumping back to the ranks of UCLA in 2013, he suffered an almost miserable start to his 2012 season,  
2063 collecting six bronze medals in 2012 and a Cambridge Olympic diploma that year. He yielded fifteen  
2064 better performances all year, which, after his limbs started teething dry, the San Francisco Nationals player  
2065 "autureted out".

2066 Analytical reports have found nothing happened, and he faces a battle to climb back on his team loaners  
2067 following his November 11th release by the Massachusetts Athletic Department.

2068 It have not been very easy. Less than one Friday ago, he went bad back and guilty of punching the face of a  
2069 wrestler. And in recent days, his approach to appeal has fared poorly. Last season, after a mid-season torn  
2070 fall, he tank-junked in an US court of law - with his lawyers threatening to sue him unless he addressed his  
2071 alleged doping confession publicly.

2072 The criticism was even louder this season.

2073 Image copyright Getty Images Image caption Jovan Penscault is in turmoil in the NBA

2074 Ask a few probing questions, in Belnod's case as well; he is asked to explain his appetite, his courtship  
2075 with his parents and whether there had been a bad day in the offseason. Yet the assumption is not that he is  
2076 emotionally disturbed. He is not unfaithful or arrogant. We believe, in his eyes, that outdresses are just trying  
2077 to achieve their own goals.

2078 Character can have positive effects, however, and even if there are further positive developments in the NBA,  
2079 the progression of such a controversial player is a critique of how US sports and the game works. It is certainly  
2080 a negative sign and if anything ought to change about how it works to bolster a program.

2081 "Coaches used to talk about their tennis stars in majors. Now they've got them everywhere," says Gary Smith,  
2082 the head coach at Penn State, who doesn't overstate the size of the reward for those who played in the nation's  
2083 top SSQC or were willing to put in athletic service. "Look at what a lot of your better players did. I mean,  
2084 you can't

2085 ===== Prompt Length 512 =====

2086 with the Hawks with a view to identifying why, as they say in the song, Hawthorn is such a happy team.  
2087 Which raises the question: Are the Hawthorn players happy because they have won the past two premierships  
2088 or is Hawthorn winning more because their players are happy? Or, as the Hawks themselves suspect, perhaps  
2089 the answer lies between. Not only does that club feel validated in their work practices, which this season have  
2090 included an unscheduled day off after their first Launceston game, but it could also use its position at the top  
2091 of the off-field ladder as it works behind the scenes to mount a case against the AFL's move in taxing clubs'  
2092 football department spending. That new tax will be reviewed before the end of next season and Hawthorn are  
2093 looking to lobby for the exclusion of such welfare initiatives such as the regular monitoring of every player's  
2094 mental health by its sport psychologist from the new tax. Another area the Hawks believe should be exempt  
2095 is the sponsoring of international educational trips for its staff. Hawthorn is not the only team that claims to  
2096 be focusing on better work-life balance for their players but they have worked more diligently to identify red  
2097 flags among their team since Travis Tuck and his mental-health issues became front-page news by way of a  
2098 third positive drug strike.

2099 While the players' union continues to investigate the reasons behind the grievances of so many of its members,  
2100 it appears beyond doubt that the demands of pre-season training have taken a disproportionate toll. This,  
2101 along with the pressure most footballers feel, even during the bulk of their holiday periods, to report back for  
2102 work in perfect condition. The inescapable conclusion drawn by Marsh and his team is that the expectations  
2103 placed upon players over their recently increased leave period have created a pre-season before the official  
2104 pre-season. According to the AFLPA and the increasingly shared belief of their leading players, the game  
2105 as a spectacle would not be adversely affected by a less-intensive spring-summer. Australian football, after  
2106 all, is a domestic sport. And the union are also monitoring, as reported by Fairfax Media, the link with the  
2107 growing injury toll. Marsh was reluctant to detail individual club player complaints but, according to his key  
2108 executive Ian Prendergast, the overall picture of dissatisfaction with their sport by the players "hit him right  
2109 between the eyes". Marsh came to the AFL from cricket and lack of enjoyment with the game they took on  
2110 through love and fun and enjoyment as children was relatively non-existent among Australian footballers  
2111 in general. Over the past decade or so, the dissatisfaction grew. They were routinely asked how much they  
2112 needed to goof around by their parents and by their coaches, which they tended to miss. Naturally, the players

2113  
2114

2115 often had their doubts about their first grade standard and few fans made allowances for this chipper criticism  
2116 when they backed their team throughout the long pre-season preparation period.

2117 The AFL and NRL clubs have recently been struggling with growing talk of missed kickouts and sudden  
2118 increments in the grading system after their clubs were caught failing to properly police Josh Boyd and  
2119 Scott Lilly. In Canberra, where some of the players are running faster and talking more about drug use,  
2120 the federation is testing the AFL's new anti-discrimination law, how to change the rules governing strip  
2121 derangement and the expanded power of post-match sanctioning cameras. Before games, the representatives  
2122 from the Sydney Maroons board, including their New Zealand captain and Sydney's assistant coach while  
2123 in charge, have delivered long diatribes to the players, to the players themselves, some of which severed  
2124 their ties to the club in return for the shortsighted and grievous attacks on the club's reputation and its unity.  
2125 There was considerable assistance from Sydney in counting the tally of the AFLU's members, from the club's  
2126 financial manager Damien Gray and chief operating officer Humphry Rokakley, to ASADA director-general  
2127 Roderic del Fuerte and former AFLAUT president Chris Jacobs. Unfortunately, Stephen Parish, who became  
2128 the Australian AFL executive president once and for all, has told a Senate inquiry that it would be "unusual  
2129 for any player not to have remitted from this federation's support when they first entered the competition  
2130 formats - and especially so if most do not do so immediately in an effort to improve their competitiveness  
2131 personally". Australian AFLU president Arthur Hurley and the AFLU's chief executive, Stephen Lawrence,  
2132 have recently left. Retired federal team officials and CEO Stephen Hendell, Hurley and Lawrence were all  
2133 also members of the most recent AFLU national advisory committee of the WFWA, the new confederation  
2134 formed before the federal government's new sporting governance regime. The refs and some of the former  
2135 clubs are quite possibly expected to have already caught on and most have said that the new local-territory  
2136 administration is very unlikely to take robust legal action against the only surviving Australian football union.  
2137 According to Jeff Pearce, the AFLPA's national vice president, the league's new federal strategy now outlines  
2138 the significant opposition to place on competitive

2138 ===== Prompt Length 768 =====

2139 with the Hawks with a view to identifying why, as they say in the song, Hawthorn is such a happy team.  
2140 Which raises the question: Are the Hawthorn players happy because they have won the past two premierships  
2141 or is Hawthorn winning more because their players are happy? Or, as the Hawks themselves suspect, perhaps  
2142 the answer lies between. Not only does that club feel validated in their work practices, which this season have  
2143 included an unscheduled day off after their first Launceston game, but it could also use its position at the top  
2144 of the off-field ladder as it works behind the scenes to mount a case against the AFL's move in taxing clubs'  
2145 football department spending. That new tax will be reviewed before the end of next season and Hawthorn are  
2146 looking to lobby for the exclusion of such welfare initiatives such as the regular monitoring of every player's  
2147 mental health by its sport psychologist from the new tax. Another area the Hawks believe should be exempt  
2148 is the sponsoring of international educational trips for its staff. Hawthorn is not the only team that claims to  
2149 be focusing on better work-life balance for their players but they have worked more diligently to identify red  
2150 flags among their team since Travis Tuck and his mental-health issues became front-page news by way of a  
2151 third positive drug strike.

2152 While the players' union continues to investigate the reasons behind the grievances of so many of its members,  
2153 it appears beyond doubt that the demands of pre-season training have taken a disproportionate toll. This,  
2154 along with the pressure most footballers feel, even during the bulk of their holiday periods, to report back for  
2155 work in perfect condition. The inescapable conclusion drawn by Marsh and his team is that the expectations  
2156 placed upon players over their recently increased leave period have created a pre-season before the official  
2157 pre-season. According to the AFLPA and the increasingly shared belief of their leading players, the game  
2158 as a spectacle would not be adversely affected by a less-intensive spring-summer. Australian football, after  
2159 all, is a domestic sport. And the union are also monitoring, as reported by Fairfax Media, the link with the  
2160 growing injury toll. Marsh was reluctant to detail individual club player complaints but, according to his key  
2161 executive Ian Prendergast, the overall picture of dissatisfaction with their sport by the players "hit him right

2162 between the eyes". Marsh came to the AFL from cricket and lack of enjoyment with the game they took on  
2163 through love and fun and enjoyment as children was relatively non-existent among Australian cricketers.  
2164  
2165 "Maybe I'm being romantic about sport, but I do think it should be fun," was all Marsh would say on the  
2166 subject. "It's a question which needs to be addressed by everyone in the industry." The AFL is not unique  
2167 in that they seems paralysed to act against the increasing pressure and negative impact of social media and  
2168 the uglier side of the expanded modern "selfie" syndrome so hauntingly articulated by Chris Judd in these  
2169 pages. But the AFLPA have identified player grievances it can address in their next wages-and-conditions  
2170 deal. Heading those grievances are the plethora of meetings the players believe over-punctuate their working  
2171 week and many believe have grown as assistant coaches over-zealously justify their jobs. Another is the  
2172 growing demand of player appearances by clubs working to service members and sponsors. The concerns  
2173 of the players' representative body were compounded by their pre-season round of visits but were already  
2174 being addressed after the second annual players' survey last August. Players responded to questions across  
2175 three areas: their club's culture; resources and its structure, which included leave periods, time off and the  
2176 performance of player development managers. The AFLPA reported back to the clubs, telling them where they  
2177 were ranked in each of those four areas. The annual research session was six-day. Those clubs in these four  
2178 areas that showed the worst results at all times and where their own poor progress was identified as the most  
2179 apparent, such as the Giants, Hawthorn and Carlton, were laid-out as operating not only without the frequent  
2180 parliamentary supervises and the more outspoken directorial redistribution policy of the Granny Hill clubs  
2181 – and where the players most often are active in their pursuit of the best and fairest – but they also operate  
2182 with the most lax end-to-end performance-related coaching policies, have the most overworked executive list  
2183 management and team performance-management staff, and rank among the worst in their management of  
2184 ticket market participation and other competition activities. Determination of future action at many clubs  
2185 were not only diluted in terms of the performance of key people in the sports knowledge-based management  
2186 structures – such as communication engineers, performance support managers and the managers of the union  
2187 candidates, but also the clubs' player development agents, tryouts and on-balling cadres. Plus, this was  
2188 punctuated by the poor performances of the club's executive list manager and model player association  
2189 executive coach – who was rated below average.

2188 "The Blues said they wanted

2191 H.2 SAMPLES GENERATED BY THE VARIABLE-LENGTH MODEL TRAINED ON STORIES

2192  
2193 ===== Sample 1 =====  
2194  
2195 once there was a boy who liked to explore - even though he was small. every day he ' d struggle to see  
2196 what he could find. he would pick small foods in his garden and visit his neighbours. one day, he noticed a  
2197 delicious sauce on the food. he was so tempted to try it out! he struggled to eat it, not as usual, so he decided  
2198 to buy a popular pasta with sauce. he shared it with all his friends who enjoyed their sauce. but, he refused  
2199 when he urged them to keep buying. in consideration, someone called out to him : " go ahead, try the sauce.  
2200 what if the sauce means you won ' t have to really savor the taste? " the boy was happy that he couldn ' t  
2201 resist the temptation to try the sauce. from then on, everyone in the neighbourhood remembered the popular  
2202 sauce and dreamed of becoming a more daring sauce expert.

2203 ===== Sample 2 =====  
2204  
2205 john was a very popular three year old boy in his town. he liked to play in his backyard in the grass. every  
2206 day, he would imagine all the different things the world started to offer. one day, john looked up in the sky  
2207 and came up with an idea. he got off his bike to learn the wisdom. he rode back to the meadow and told the  
2208 wise owl about the wisdom. the owl said she said the world was full of special and sweet wisdom if he was  
patient. after a long day, the wisdom yielded a lot. john thanked the wise owl. she had worked hard and was

2209 able to learn something magical. john continued on his bike, soon out of the meadow. he saw a beautiful  
2210 butterfly and smiled as attractive as she had seen him.  
2211

2212 ===== Sample 3 =====

2213 tim was playing in the park one day. he skipped around, shouting. it was a game but there were no toys  
2214 around. he was feeling disappointed and soon saw a boy. the boy was a lot bigger than no other boy but no  
2215 one ever seemed to understand. he recognized it from a playground. the boy asked if he wanted to join him.  
2216 tim hesitated at first for a few moments. but he felt a bit scared, so he decided to follow the boy. he started to  
2217 climb up the ladder at first but as he went higher and higher he felt helpless. the boy saw tim and said he was  
2218 looking for help. tim was gentle and asked to try to lift him up the ladder but he was very worried. after a  
2219 while the boy said goodbye to tim and ran off. tim felt sad that he had to go home without to help a friend. he  
2220 regrets climbing the ladder.

2221 ===== Sample 4 =====

2222 timmy and lily are friends. they like to play together. one day, they see a lady on a bench. she has a big bag  
2223 and a red stick. the candy is sweet and bitter. lily wants to eat candy too. she reaches into her bag and takes  
2224 the candy. " no, lily! " timmy says. " that is my candy. you can ' t have it. i want this candy. it is a candy for  
2225 you. " lily did not listen. she does not want to share her candy. she tries to part with her candy with the stick.  
2226 but the stick is stuck. it stays loose and hurts her teeth. " ow! " lily cries. " you broke my stick! " timmy  
2227 laughs. he drops the stick and pushes lily away. " no, thank you! " lily says. " you can ' t have my candy. go  
2228 away now. " tom does not find another stick. he did not know how to have it. he bites the stick and pulls it  
2229 out. " look, lily, i am eating candy! " tom says. " am i eating kite? " lily looks at tom. she saw what tom was  
2230 doing. she feels silly and scared. she takes off her hat and her shoes and pants. " oops! sorry, tom, " timmy  
2231 says. " i did not mean to scare you. i just wanted to bite the stick. maybe we can play something else. " " i  
2232 know, tom, " lily says. " but i still like the stick. it is better. you can have another candy. " tom feels bad. he  
2233 also likes lily. he thinks lily is generous. he gave lily a big, red apple. " here, lily, " he says. " this is a cake. it  
2234 was not magic. it was a secret. " " thank you, lily, " tom says. " but next time, do not eat the stick. you are a  
2235 good friend. can i help you? " lily smiles. she is glad that tom could help. she was generous. " okay, lily, "  
2236 tom says. " i will share the stick. we can play together. " they play with a ball, a kite, a book and other toys.  
2237 they have fun with the ball. they are happy.

2238 ===== Sample 5 =====

2239 once upon a time there was a boy named jack. he was worried because he had no worries with him, every day  
2240 him would carry jack away for a day at work the first time to eat a yummy lunch. one day jack saw some  
2241 of his mum outside. they were so kind and helpful to all of them. they did not put a heavy box at school  
2242 or they brought any more delicious treats. jack told his mum about this because he went to work looking  
2243 for something else to do. jack soon had an idea : he bought a lunchbox and place all of his mum ' s lunch.  
2244 he carried each piece, and then drive home. his mum was so relieved and smiled. jack was so relieved and  
2245 happy. from that day on, jack always carried his lunch with him everyday, and people never forgot to go home  
2246 without a worried look.

2247 ===== Sample 6 =====

2248 once there was a dog named spot. spot had a big toy that was very hairy. he loved to play with it and hug him.  
2249 every day spot would sign when he would get a slapped race. one time, spot was ready to show his signa€  
2250 " black pawch. he wanted to show his friends how fast he could sign? when he signed with his hand, his  
2251 friends saw him close by. they liked the show. the show was very funny. when spot managed to sign his name,  
2252 his friends would even make clapping. in the end, spot and his friends all got slap spots! they are now the  
2253 best! all day long, spot is signing, and even his hand feels good. everyone around him is happy every week.  
2254 because of the day, spot brought a prize for his slapch with him.  
2255



2256 ===== Sample 7 =====  
2257  
2258 once upon a time there was a gorilla. he was big and very strong. he stored food in a crack. the gorilla enjoyed  
2259 the food and enjoyed it. one day, a little boy was watching the gorilla by himself. he got stuck in a crack. he  
2260 felt scared and demanded the gorilla some help from it. he knew what to do. he asked the penguin for help.  
2261 the penguin grabbed a big rock and swam over to the crack. the gorilla worked very hard with his help from  
2262 the penguin. the little boy was so happy. he ate his food and the penguin thanked him for the sweet treat.

2263 ===== Sample 8 =====  
2264  
2265 once upon a time, there was a little girl named lily who loved to play with her ball. every day, she would go  
2266 to the park and send her ball to him inside a loud song. one day, lily saw a mean boy named max playing  
2267 his guitar. he was playing his song and was not friendly. lily wanted to help max, so she played with him  
2268 and made the music beautiful. max became friendly and welcomed lily with her ball. they had a great time  
2269 together and at the end of the musical song, max ended nicely and said goodbye. lily felt sad that he was not  
2270 sharing nicely. lily went home and used a plan. she brought out her guitar and spoke to max. she told him  
2271 that he wasn't nice and allowed her to hold the guitar. max was happy and grateful. from that day on, lily  
2272 and max didn't matter what friends tried from each other.

2273 ===== Sample 9 =====  
2274  
2275 once upon a time, there was a messy little boy named timmy. timmy liked to play with his toys and do chores.  
2276 one day, timmy's stomach felt really hungry. his tummy was growling, and his food was all over the floor.  
2277 his mom saw that timmy had eaten all his food, and she knew the bad smell came from lunch. timmy ate it,  
2278 but she found it was disgusting and looked at the napkins on the floor. she said, " timmy, napkins are not good  
2279 to eat. " she scolded timmy and asked him again before eating his meals. she said, " you don't know you  
2280 would have gotten sick, though. " timmy said, " i'm sorry, i won't eat them next time. " he put his hands on  
2281 the napkins and threw them away so his clothes were safe from harm. from then on, timmy made sure to listen  
2282 to his mom and eat his napkins before she went to eating. he never wanted to eat his napkins again. the end.

2283 ===== Sample 10 =====  
2284  
2285 anna liked to play with the sack. she liked to pull things out and see what she could find. she asked her dad  
2286 sometimes. he explained that she was anything she could find. one day, anna found the sack in the field. she  
2287 looked inside and saw her friend, lily, who was playing with a dog. she opened the sack and saw many shiny  
2288 things. anna was happy. she wanted to explore more. but then, she saw a big black cat. the cat was sitting on  
2289 the edge of a bush. it had dirty eyes. anna did not know what it was. she thought, " maybe it is the cat. i can  
2290 make it happy. " she went to the shiny thing and stroked it gently. the cat did not go away. it hissed. it said, "  
2291 no, anna, what are you doing? go out! let's stay behind the tree. " anna did not listen. she stepped closer to  
2292 the cat. she reached her hand out. she pulled its head. the shiny thing moved. it made a ring. it came alive  
2293 with a band. it moved and snapped. anna heard the snap. she screamed. she ran back. the cat bit her. it hurt  
2294 her finger. she could not hold it. she cried, " ow, ow, ow, ow, it hurts! it hurts, it hurts! " anna ran to the field.  
2295 she saw the fence and the other people. she saw the cat too. she heard the people shouting and running. she  
2296 saw lily's finger. she felt sorry for her. she gave her ring to her. she went to lily and said, " i'm sorry, mom. i  
2297 was wrong. lily was naughty. she tried to grab lily's ring. " her mom smiled and hugged lily. she said, " it's  
2298 okay, anna, it's okay. it doesn't hurt. but it's not your fault. you should have left the sack alone. " anna  
2299 was still sad. she said, " i'm sorry, mom. i love you, lily. " she hugged lily and learned her lesson. she never  
2300 came back. she hugged lily and went home. she left the cat alone. where lily was sad. she did not harm it.

2301 ===== Sample 11 =====  
2302  
2303 once upon a time, there was a little girl named lily. she loved to play with her toys and pet dog, max. one day,  
2304 she went to the park with her mom and saw a leopard flying from its back to its home. she excitedly asked her

2303 mom, " what is that, mommy? " her mom said to lily, " this is polly. this leopard is very cheap, so if it misses  
2304 its home, i can take it to you. " lily listened carefully and said, " thank you, mommy. polly is so kind. can i  
2305 recommend it to a party for you at the park? " later that week, lily and max watched the leopard and when  
2306 it arrived, the leopard landed on lily ' s fur. she was so excited! she learned that the leopard had a modest  
2307 quality on its home. later that day, lily forgot about the leopard and went to bed for a nap. she fell asleep with  
2308 max, just like the modest parrot. when she woke up, she thanked max for helping her.

2309 ===== Sample 12 =====

2310 once upon a time there was an ancient, grey kangaroo. he lived in a sunny spot near a pond. one night he  
2311 heard a magical sound. it was a voice coming from the sky. startled, the kangaroo called out. he didn ' t know  
2312 what to do. but then he noticed a small, red balloon glimmering in the air. he curiously approached the alien  
2313 and it turned into a truck blocking his direction. it was carrying an old gas truck and it was perfectly heavy. "  
2314 what is this? " the voice replied. the gasoline truck quickly zaneded the balloon and released it into the sky.  
2315 the kangaroo almost entered an ancient castle and it was even more beautiful than than he had ever been there  
2316 before. suddenly the gas truck spoke with a hop and the voice said, " this is a gas castle! you can ' t come  
2317 back! " the kangaroo was very scared, but he had to figure out the alien ' s names. he ran higher and higher,  
2318 and the dragon slowly started to slowly spring until it was gone. the kangaroo thought it was gone. the next  
2319 morning, the ancient kangaroo returned to where the alien had returned. he looked back, not clear with fear -  
2320 he was even more relieved to see that the magical creature had given him back. he had emerged in the right  
2321 circle!

2322 ===== Sample 13 =====

2323 once upon a time, there was a little rabbit named benny. benny loved to hop around the forest and play with  
2324 his friends. one day, benny met a friendly rabbit named rosie. benny introduced rosie to his new friend. "  
2325 what ' s that? " he asked. " that ' s a toy staff, " said rosie, " it ' s a loud sound that makes me reverse. " benny  
2326 didn ' t understand what that meant, but rosie continued. benny asked rosie if he wanted to teach rosie her  
2327 name staff. rosie said, " i don ' t know it, a toy staff means my number 10.... " benny smiled and said, " okay,  
2328 let ' s count to ten. " benny and rosie learned how to reverse from ten over half cards. they had so much fun  
2329 counting and learning, they played with the toy staff all day long. when benny and rosie grew up playing,  
2330 they had so much fun together and were able to play with their new name toy staff.

2331 ===== Sample 14 =====

2332 once upon a time there was a little boy named john. john had a grain wheat mill, who worked there. one day  
2333 john went to the mill to make flour. he picked out some flour and put the wheat in the corn mill. as he ran  
2334 over to his house, the flour slowly churned and added it to a slice of the mill. john smiled at what he had done.  
2335 when he was finished when he rolled out his door, he met another boy called bob. bob said, " would you want  
2336 to play with me? ". john smiled and said, " yes! " so they had fun playing and running around the mill. when  
2337 john was done, bob knew exactly how bob had done it. he held up the grain and carried them both hands.  
2338 he said, " we have to take care of our wheat and be balanced for what we need ". john was balanced with  
2339 excitement. he welcomed him, gave him a pat on the shoulder of my mill and was extremely pleased. bob  
2340 laughed, and told john that day he and the wheat were friends forever.

2341 ===== Sample 15 =====

2342 ben and mia like to play on the navy boat on the sea. they pretend they are divers and fight sharks in space.  
2343 one day, they go to a big island with their house. it is very pretty and green. it has a long sail, a flag, and a  
2344 blue door. they open the door and see a lot of water in the water. they lower the boat and look over the boat. "  
2345 hello, many ships! " mia says. " maybe they are pirate? " ben wonders. " maybe they are treasure, " mia says.  
2346 " or other ships. we are swimming and looking for something to look for. " they see a big ship in the water. it  
2347 is red and yellow and has a sword. " maybe it ' s a land from here, " ben says. " what ' s a port? " mia says. "  
2348 maybe it ' s a pirate! " " who is the treasure? " ben says. he dived close and lands on the ship at mia. he lands  
2349

2350 on mia next and puts on a scarf. she smiles and laughs. " did you hear the sound of the treasure? " ben asks. "  
2351 one, two, 3, blast off! " mia says. she smiles and jumps high. she sees the sky, the sun, a shark, a butterfly, a  
2352 shark, and a duck. they jump off the ship and run to the shore. they call their mom and dad. " mom, look what  
2353 we found! " ben says. " we are in space! " mia says. she hugs their mom and dad. they tell their parents about  
2354 their adventure. " we dive in the air and dive with a ship, " ben says. " we are great and brave explorers! "

2355 ===== Sample 16 =====

2356  
2357 once upon a time there was a young boy named sam. sam was very happy because he had a nice bench. every  
2358 day he used it to sit in the park. one day, he was sitting in the park on the bench while a boy came and slapped  
2359 his hand. " ouch " the boy said the boy. sam looked up in surprise and started to cry. the kind boy said, " i  
2360 have an idea to help the young boy. you can use your cane to make hi! it ' s all okay ". sam wiped up his tears  
2361 and looked at the young boy in the bench. it was a kind surprise that the boy had seen him before. he said, "  
2362 you ' re not as good as nice as the bench you like and why did who slap my hand? ". the young boy smiled  
2363 and told sam that he quarrel with a friend. he said, " let ' s write down my bench and draw it. it looks  
2364 like a special book ". the young boy smiled and said, " yes, let ' s draw on the bench ". sam and his friend  
2365 spent lots of time playing and talking. they wrote and drew the bench and gave it a hug. sam and the boy  
2366 enjoyed the bench together for many years went by and their canes were over, so they were happy.

2367 ===== Sample 17 =====

2368 once there were two friends, daniel and norman. they were mighty friends, who loved to play baseball. they  
2369 would practice their best in every golf game. even on this competitive day, norman still wanted to score a goal  
2370 without not being four goals as rough as norman. one day, daniel sneaked into norman ' s lead up getting his  
2371 best grades. johnnie said, " let ' s invite our team to try it out. in the starting game! " the team was excited,  
2372 they called the veterinarian. norman dashed to the chess room and daniel asked for a tough stick. johnnie took  
2373 out a tough stick and finally swung it against his friend in the second seat. they both cheered as norman was  
2374 proud that they could continue on their physical soccer game. at last, they were able to score successfully.

2375 ===== Sample 18 =====

2376 once upon a time, there was a little boy named timmy. timmy had a toy weapon, a brown sword. one day,  
2377 timmy ' s mom came in and asked him to stop playing with his sword. but timmy didn ' t let go of his weapon  
2378 and wanted to show her where he found it. " timmy, your weapon is on your bed, " his mom said mad, but  
2379 timmy didn ' t listen. later, timmy ' s mom asked him to clean up my room and put it away. " mom said i will  
2380 shoot away from your car, " she replied. timmy did a great job cleaning up the room with his dolls and his  
2381 weapon was very messy. but when he came back, his sister came in and saw that all of his toys were gone.  
2382 timmy couldn ' t find his toy anywhere. he looked in the closet, in the closet, and even in the closet. his sister  
2383 searched everywhere, but she couldn ' t find it. timmy felt sad because he didn ' t want to make her worry.  
2384 finally, he told his sister that he did not think he was stupid of any extra toys. his sister decided to find his toy  
2385 and they searched the treasure together. the sword was found at its spot and timmy was glad no one could  
2386 spot it.

2387 ===== Sample 19 =====

2388 once upon a time there was a modern woman. she had a very important job and she was very careful with it.  
2389 each day she looked at her job, as it was a great job. then one day, the woman heard a loud noise. she looked  
2390 around, looking for help. she started to look out there. suddenly, she saw a little girl running towards a 3 year  
2391 old policeman. she was running and rushing to avoid her. the policeman stopped and said, " are you alright? i  
2392 ' m so scared and lost. " the elderly woman went to find the cop and saw how brave the little girl was. he  
2393 asked her if she wanted to help. the little girl told her yes and the cop started hunting. but the little girl bit her  
2394 and threw everything away. in the end, the woman ' s heart was left increased and into distress. the moral of  
2395 the story is that we must always be careful when we are hunting, especially when someone is about three  
2396 years old and the little girl needed help. if we argued with someone, or else remains in our life, we can start

2397 causing consequences. this story requires a serious ending, figuring out those who are looking out for help  
2398 and can cause trouble.

2399  
2400 ===== Sample 20 =====

2401 once upon a time, lily and max were best friends. one day, they found a delicate flower in the grass. it was  
2402 very delicate with lots of spiky petals. max loved flowers and wanted to hold it. but the flower was too big for  
2403 him and max couldn ' t hold it very much. but then, they found a snapwig and the flower fell and landed in  
2404 their hands. it burst in two! lily was so happy that max found a beautiful flower they could rectangle. they  
2405 decided to build a castle in the castle. lily looked for pretty stones as they found, but even though it was hard  
2406 for them to find that stones would be making the castle stronger. so, they promised to keep the castle strong  
2407 and delicate from one day on. they played together every day, and from that day on, the delicate flower made  
2408 its castle strong and attractive.

2409  
2410 H.3 DEMONSTRATION OF THE INTERMEDIATE GENERATION PROCESS

2411  
2412 ===== Step 0 =====

2413 [CLS]

2414  
2415 ===== Step 1 =====

2416 [CLS] lucy

2417  
2418 ===== Step 2 =====

2419 [CLS], lucy her

2420  
2421 ===== Step 3 =====

2422 [CLS], lucy very her.

2423  
2424 ===== Step 4 =====

2425 [CLS], lucy lucy very her.

2426  
2427 ===== Step 5 =====

2428 [CLS], lucy smile lucy, very her.

2429  
2430 ===== Step 6 =====

2431 [CLS], lucy. smile lucy she, very her.

2432  
2433 ===== Step 7 =====

2434 [CLS], lucy. smile lucy she, very her.

2435  
2436 ===== Step 8 =====

2437 [CLS], lucy. smile. lucy she and, very her.

2438  
2439 ===== Step 9 =====

2440 [CLS], lucy. it smile. lucy she and, very her.

2441  
2442 ===== Step 10 =====

2443 [CLS], lucy. it smile. lucy she and, very after her.

2444  
2445 ===== Step 11 =====

2444 [CLS], lucy. it smile. lucy she and, very after her. her  
2445 ===== Step 12 =====  
2446 [CLS], lucy.. it smile. lucy she and, very after her. park her  
2447 ===== Step 13 =====  
2448 [CLS], lucy.. it smile. lucy she and, very after her. park her.  
2449 ===== Step 14 =====  
2450 [CLS], lucy. a. it smile. lucy she and, very pretty after her. park her.  
2451 ===== Step 15 =====  
2452 [CLS], lucy. a. it smile. lucy she and farm, very pretty after her. park her.  
2453 ===== Step 16 =====  
2454 [CLS], lucy. a. it smile. lucy she the and farm farm, very pretty all after finished her. park her.  
2455 ===== Step 17 =====  
2456 [CLS], lucy. a. it smile. lucy to she the and farm farm,. very pretty playing all after finished her. park her.  
2457 ===== Step 18 =====  
2458 [CLS], there lucy. a. it smile. lucy to she the and farm farm,. very pretty playing all after finished her. park  
2459 her.  
2460 ===== Step 19 =====  
2461 [CLS], there lucy. a. it smile. lucy to she the and farm farm,. onion very pretty playing all after finished her.  
2462 to park her.  
2463 ===== Step 20 =====  
2464 [CLS], there lucy. a. it and smile. lucy to. she the and farm farm,. onion very pretty playing all after finished  
2465 her. to park her mom.  
2466 ===== Step 21 =====  
2467 [CLS], there lucy. a. it and smile. lucy to. she the and farm farm,. onion very pretty playing all after finished  
2468 her. to park her mom more.  
2469 ===== Step 22 =====  
2470 [CLS], there lucy. a. it and smile. lucy to. she saw the and farm farm,. onion very pretty playing all after  
2471 finished her. to park her mom the more fun.  
2472 ===== Step 23 =====  
2473 [CLS] a, there lucy. a. it and smile. lucy the to. she saw the and farm farm,. onion very pretty playing all after  
2474 finished her. to park her mom the more fun.  
2475 ===== Step 24 =====  
2476 [CLS] a, there lucy. a. it and smile. lucy the to. she saw the and farm farm,. onion very pretty playing all after  
2477 finished her. to park her mom the more fun.  
2478 ===== Step 25 =====  
2479  
2480  
2481  
2482  
2483  
2484  
2485  
2486  
2487  
2488  
2489  
2490

2491 [CLS] a, there lucy. a. it and smile. lucy the to. she saw the and farm farm,. onion very pretty playing all after  
2492 finished her. to park her mom the more fun.  
2493  
2494 ===== Step 26 =====  
2495 [CLS] a, there lucy. a. it and smile. lucy the to. she saw the and farm farm,. onion very pretty playing all after  
2496 finished her. to park her mom the more fun.  
2497  
2498 ===== Step 27 =====  
2499 [CLS] a, there lucy. a. it and smile. lucy the park to. she saw the and farm farm,. onion very pretty playing all  
2500 after finished, her. to park her mom the they more fun.  
2501  
2502 ===== Step 28 =====  
2503 [CLS] a, there lucy. a. it and smile. lucy the park to. she saw the and farm farm,. onion very pretty playing all  
2504 after finished, her. to park her mom the they more fun.  
2505  
2506 ===== Step 29 =====  
2507 [CLS] a, there was lucy. she a. it and smile. lucy the park to. she saw the and farm farm,. onion very pretty  
2508 playing all after finished, her. to park with her mom the they more fun.  
2509  
2510 ===== Step 30 =====  
2511 [CLS] a, there was lucy. she a. it attractive and smile. lucy the park to. she saw the and farm farm,. onion  
2512 very pretty playing all after finished, her. to park with her mom the they more fun.  
2513  
2514 ===== Step 31 =====  
2515 [CLS] a, there was lucy. she a. it attractive and smile. lucy the park to. she saw the and farm farm,. onion  
2516 very pretty playing all after finished, her. very to park with her mom the they more fun the.  
2517  
2518 ===== Step 32 =====  
2519 [CLS] a, there was lucy. she a. it attractive and smile. lucy the park to. she saw the and farm farm,. onion  
2520 very pretty playing all after finished, her. very to park with her mom the they more fun the.  
2521  
2522 ===== Step 33 =====  
2523 [CLS] a, there was lucy. she a. it attractive and smile. lucy the park to. she saw the and farm farm,. onion  
2524 very pretty. playing all after finished, her. very to park with her mom the they more fun the.  
2525  
2526 ===== Step 34 =====  
2527 [CLS] a, there was lucy. she a. it attractive and smile. lucy the park to. she saw the and farm farm,. onion  
2528 very pretty. playing all after finished, her. was very to the park with her mom the they more fun the.  
2529  
2530 ===== Step 35 =====  
2531 [CLS] a time, there was lucy. she a. it was attractive and smile. lucy the park to. she saw the trees and farm  
2532 farm,. onion very pretty. playing all after finished, her. was very to have the park with her mom the they more  
2533 fun the.  
2534  
2535 ===== Step 36 =====  
2536 [CLS] a time, there was lucy. she a. it was attractive and smile. lucy the park to. she saw the trees and farm  
2537 farm,. onion very pretty. playing all after finished, her. was very to have the park with her mom the they more  
fun the.

2538 [CLS] a time, there was lucy. she a. it was attractive and smile. lucy the park to. she saw flowers the trees and  
2539 the farm farm,. onion very pretty. playing all after finished, her. was very to have the park with her mom the  
2540 they more fun the.

2541 ===== Step 38 =====

2543 [CLS] a time, there was lucy. she a. it was attractive and smile. day lucy the park to. she saw flowers the trees  
2544 and the farm farm, onions. onion very pretty. playing all after finished, went her. was very to have the park  
2545 with her mom the onion they more fun the.

2546 ===== Step 39 =====

2547 [CLS] a time, there was lucy. she a. it was attractive and made smile. day lucy the park to. she saw flowers  
2548 the trees and the farm farm, onions. onion very pretty. playing all after finished, went her. was very to have  
2549 the park with her mom the onion they more fun the.

2551 ===== Step 40 =====

2552 [CLS] a time, there was lucy. she a. it was very attractive and made smile. day lucy the park to. she saw  
2553 flowers the trees and the farm farm, onions. onion very pretty. playing all after finished, went her. was very to  
2554 have the park with her mom the onion they more fun the.

2556 ===== Step 41 =====

2557 [CLS] a time, there was named lucy. she a. it was very attractive and made smile. day lucy the park to. she  
2558 saw flowers, the trees and the farm farm, onions. onion very pretty. playing all after finished, went her. was  
2559 very to have the park with her mom the onion they more fun the.

2560 ===== Step 42 =====

2561 [CLS] a time, there was named lucy. she a. it was very attractive and made smile. day lucy the park to. she  
2562 saw flowers, the trees, and the farm farm, onions. onion very pretty. playing, all after finished, went her. was  
2563 very to have the park with her mom the onion they more fun at the.

2565 ===== Step 43 =====

2566 [CLS] a time, there was named lucy. she a. it was very attractive and made smile. day lucy the park to. she  
2567 saw flowers, the trees, and the farm farm, onions. tomatoes onion very pretty. playing, all after finished, went  
2568 her. was very to have the park with her mom the onion they planned more fun at the.

2570 ===== Step 44 =====

2571 [CLS] a time, there was named lucy. she a. it was very attractive and made smile. day lucy went the park  
2572 to. she saw flowers, the trees, and the farm farm, onions. tomatoes onion were very pretty. playing, all after  
2573 finished, went her mom. she was very to have the park with her mom the onion they planned more fun at the.

2574 ===== Step 45 =====

2575 [CLS] a time, there was named lucy. she a. it was very attractive and made smile. day lucy went the park  
2576 to. she saw flowers, the trees, and the farm farm, onions. orange tomatoes onion were very pretty. playing,  
2577 watered all after finished, went her mom. she was very to have the park with her mom the onion they planned  
2578 more fun at the.

2580 ===== Step 46 =====

2581 [CLS] a time, there was named lucy. she a. it was very attractive and made smile. day lucy went the park  
2582 to. she saw flowers, the trees, and the farm farm, onions. orange tomatoes onion were very pretty. playing,

2583  
2584

2585 watered all after finished, went her mom. she was very to have the park with her mom the onion they planned  
2586 more fun at the.

2587  
2588 ===== Step 47 =====

2589 [CLS] a time, there was named lucy. she a. it was very attractive and made smile. day lucy went the park to.  
2590 she saw flowers, the trees, and the farm farm, onions. orange tomatoes and onion were very pretty. playing,  
2591 watered all after finished planting, went her mom. she was very to have seen the park with her mom and the  
2592 onion they planned more fun at the.

2593  
2594 ===== Step 48 =====

2595 [CLS] a time, there was named lucy. she a. it was very attractive and made smile. day lucy went the park to.  
2596 she saw flowers, the trees, and the farm farm, onions. orange tomatoes and onion were very pretty. playing,  
2597 watered all after finished planting, went to her mom. she was very to have seen the park with her mom and  
2598 the onion they planned more fun at the.

2599  
2600 ===== Step 49 =====

2601 [CLS] a time, there was named lucy. she a dress. it was very attractive and made smile. day lucy went the  
2602 park to. she saw flowers, the trees, and the farm farm, lucy onions. they orange tomatoes, and onion were  
2603 very pretty. playing, watered all after finished planting, went to her mom. she was very to have seen the park  
2604 with her mom and the onion they planned more fun at the.

2605  
2606 ===== Step 50 =====

2607 [CLS] a time, there was named lucy. she a dress. it was very attractive and made smile. day lucy went the  
2608 park to. she saw flowers, the trees, and the farm farm, lucy onions. they orange tomatoes, and onion were  
2609 very pretty. playing, watered all after finished planting, went to her mom. she was very to have seen the park  
2610 with her mom and the onion they planned more fun at the.

2611  
2612 ===== Step 51 =====

2613 [CLS] once a time, there was named lucy. she a dress. it was very attractive and made smile. day lucy went to  
2614 the park to. she saw flowers, the trees, and the farm farm, lucy saw onions. they orange tomatoes, and onion.  
2615 were very pretty. playing, watered all after finished planting, went to her mom. she was very to have seen the  
2616 park with her mom and the onion they planned more fun at the day.

2617  
2618 ===== Step 52 =====

2619 [CLS] once a time, there was named lucy. she a dress. it was very attractive and made her smile. day lucy  
2620 went to the park to. she saw flowers, the trees, and the farm farm, lucy saw onions. they orange tomatoes, and  
2621 onion. thought were very pretty. playing, watered all after lucy finished planting, went to her mom. she was  
2622 very to have seen the park with her mom and the onion they planned more fun at the day.

2623  
2624 ===== Step 53 =====

2625 [CLS] once a time, there was named lucy. she a dress. it was very attractive and made her smile. day lucy  
2626 went to the park to. she saw flowers, the trees, and the farm farm, lucy saw onions. they orange tomatoes, and  
2627 onion. thought were very pretty. playing, watered all vegetables after lucy finished planting, went to her mom.  
2628 she was very happy to have seen the park with her mom and the onion they planned more fun at the day.

2629  
2630 ===== Step 54 =====

2631 [CLS] once a time, there was little named lucy. she a dress. it was very attractive and made her smile. day  
lucy went to the park to. she saw flowers, the trees, and the farm farm, lucy saw onions. they orange tomatoes,  
and onion. thought were very pretty. after playing, watered all vegetables after lucy finished planting, went to



2632 her mom. she was very happy to have seen the park with her mom and the onion they planned more fun at the  
2633 every day.

2634  
2635 ===== Step 55 =====

2636 [CLS] once a time, there was little named lucy. she a dress. it was very attractive and made her smile. day,  
2637 lucy went to the park to. she saw flowers, the trees, and the farm farm, lucy saw onions. they orange tomatoes,  
2638 and onion. thought were very pretty. after playing, lucy watered all her vegetables after lucy finished planting,  
2639 went to her mom. she was very happy to have seen the park with her mom and the onion they planned more  
2640 fun at the park every day.

2641  
2642 ===== Step 56 =====

2643 [CLS] once a time, there was little named lucy. she a dress. it was very attractive and made her smile. one  
2644 day, lucy went to the park to. she saw flowers, the trees, and the farm farm, lucy saw onions. they orange  
2645 tomatoes, and onion. thought they were very pretty. after playing, lucy watered all her vegetables after lucy  
2646 finished planting, went to her mom. she was very happy to have seen the park with her mom and the onion  
2647 they planned more fun at the park every day.

2648  
2649 ===== Step 57 =====

2650 [CLS] once a time, there was a little named lucy. she a nice dress. it was very attractive and made her smile.  
2651 one day, lucy went to the park to. she saw flowers, the trees, and the farm the farm, lucy saw onions. they  
2652 orange tomatoes, and onion. thought they were very pretty. after playing, lucy watered all her vegetables after  
2653 lucy finished planting, went to her mom. she was very happy to have seen the park with her mom and the  
2654 onion they planned more fun at the park every day.

2655  
2656 ===== Step 58 =====

2657 [CLS] once a time, there was a little named lucy. she a nice dress. it was very attractive and made her smile.  
2658 one day, lucy went to the park to. she saw the flowers, the trees, and the farm the farm, lucy saw onions. they  
2659 orange tomatoes, and onion. thought they were very pretty. after playing, lucy watered all her vegetables after  
2660 lucy finished planting, went to her mom. she was very happy to have seen the park with her mom and the  
2661 onion they planned more fun at the park every day.

2662  
2663 ===== Step 59 =====

2664 [CLS] once upon a time, there was a little girl named lucy. she had a nice dress. it was very attractive and  
2665 made her smile. one day, lucy went to the park to play. she saw the flowers, the trees, and the farm. the farm,  
2666 lucy saw many onions. they orange tomatoes, and onion. thought they were very pretty. after playing, lucy  
2667 watered all her vegetables after lucy finished planting, went to her mom. she was very happy to have seen the  
2668 park with her mom and the onion they planned more fun days at the park every day.

2669  
2670 ===== Step 60 =====

2671 [CLS] once upon a time, there was a little girl named lucy. she had a nice dress. it was very attractive and  
2672 made her smile. one day, lucy went to the park to play. she saw the flowers, the trees, and the farm. the farm,  
2673 lucy saw many onions. they orange tomatoes, and onion. she thought they were very pretty. after playing,  
2674 lucy watered all her vegetables after lucy finished planting, went to her mom. she was very happy to have  
2675 seen the park with her mom and the onion they planned more fun days at the park every day.

2676  
2677 ===== Step 61 =====

2678 [CLS] once upon a time, there was a little girl named lucy. she had a nice dress. it was very attractive and  
made her smile. one day, lucy went to the park to play. she saw the flowers, the trees, and the farm. at the  
farm, lucy saw many onions. they orange tomatoes, and orange onion. she thought they were very pretty.

2679 after playing, lucy watered all her vegetables after lucy finished planting, went to her mom. she was very  
2680 happy to have seen the park with her mom and the onion. they planned more fun days at the park every day.  
2681

2682 ===== Step 62 =====

2683 [CLS] once upon a time, there was a little girl named lucy. she had a nice dress. it was very attractive and  
2684 made her smile. one day, lucy went to the park to play. she saw the flowers, the trees, and the farm. at the  
2685 farm, lucy saw many onions. they were orange tomatoes, and orange onion. she thought they were very pretty.  
2686 after playing, lucy watered all her vegetables after lucy finished planting, went to her mom. she was very  
2687 happy to have seen the park with her mom and the onion. they planned more fun days at the park every day.

2688 ===== Step 63 =====

2689 [CLS] once upon a time, there was a little girl named lucy. she had a nice dress. it was very attractive and  
2690 made her smile. one day, lucy went to the park to play. she saw the flowers, the trees, and the farm. at the  
2691 farm, lucy saw many onions. they were orange tomatoes, and orange onion. she thought they were very pretty.  
2692 after playing, lucy watered all her vegetables. after lucy finished planting, went to her mom. she was very  
2693 happy to have seen the park with her mom and the onion. they planned more fun days at the park every day.  
2694

2695 ===== Step 64 =====

2696 [CLS] once upon a time, there was a little girl named lucy. she had a nice dress. it was very attractive and  
2697 made her smile. one day, lucy went to the park to play. she saw the flowers, the trees, and the farm. at the  
2698 farm, lucy saw many onions. they were orange tomatoes, and orange onion. she thought they were very pretty.  
2699 after playing, lucy watered all her vegetables. after lucy finished planting, she went to her mom. she was very  
2700 happy to have seen the park with her mom and the onion. they planned more fun days at the park every day.  
2701

2702  
2703  
2704  
2705  
2706  
2707  
2708  
2709  
2710  
2711  
2712  
2713  
2714  
2715  
2716  
2717  
2718  
2719  
2720  
2721  
2722  
2723  
2724  
2725