

Diffusion Language Models Generation Can Be Halted Early

Anonymous ACL submission

Abstract

Diffusion Language models (DLMs) are a promising avenue for text generation due to their practical properties on tractable controllable generation. They also have the advantage of not having to predict text autoregressively. However, despite these notable features, DLMs have not yet reached the performance levels of their Autoregressive counterparts. One of the ways to reduce the performance gap between these two types of language models is to speed up the generation of DLMs. Therefore, we propose a pioneering methodology to address this issue in this work. It enables the execution of more generation steps within a given time frame, potentially leading to higher-quality outputs. Specifically, our methods estimate DLMs completeness of text generation and allow adaptive halting of the generation process. We test and refine our methods on Plaid, SSD, and CDCD DLMs and create a cohesive perspective on their generation workflows. Finally, we confirm that our methods allow halting Plaid, SSD, and CDCD models and decrease the generation time by 10-40% without a drop in the quality of model samples.

1 Introduction

Language Models (LMs) are essential Natural Language Processing (NLP) tools. The two primary methods of training LMs for NLP are autoregressive training (Radford et al., 2019; Raffel et al., 2020; Chowdhery et al., 2022) and masked language modeling (Devlin et al., 2019; He et al., 2020; Liu et al., 2019; Lan et al., 2020).

The exploration of alternative models, such as Diffusion Models (Ho et al., 2020; Song et al., 2020), is a promising avenue for research as diffusion allows native non-causal conditioning and simplified controllable generation methods (Nichol et al., 2022). In recent works, with models such as "Diffusion LM" and Plaid (Li et al., 2022; Gulrajani and Hashimoto, 2023), Simplex-based Diffusion

Language Model (SSD) (Han et al., 2023), GENIE (Lin et al., 2022), and Continuous Diffusion for Categorical Data (CDCD) (Dieleman et al., 2022) being introduced, we can see an emerging interest for using Diffusion Models in text generation.

A crucial distinction between Autoregressive LMs and Diffusion Language Models (DLMs) lies in their modeling approaches. Autoregressive LMs predominantly adhere to the common probabilistic model. In contrast, DLMs exhibit substantial divergence in their application for modeling categorical data. When exploring DLMs, it is essential to consider the lack of connectivity between such models. The majority of comparisons between them have primarily focused on evaluating sample quality.

While it is essential to study the sample quality of DLMs, it does not further our understanding of the differences between these models. This work addresses this issue and evaluates various DLMs with a unified view of their generation process. Given this unified view, we study the dynamics of the generation process within different DLMs and focus on the changes in the samples during that process.

The main contributions of this paper can be summarized as follows:

- We showed that the generation process of most DLMs for general text generation can be halted, which makes it possible to implement an early, faster sample generation without compromising quality.
- To the best of our knowledge, we were the first to evaluate DLMs with adaptive Early Exiting (Graves, 2016). In this paper, we introduced three adaptive criteria inspired by the ones used for text classification (Liu et al., 2020; Zhou et al., 2020; Gao et al., 2023).
- We evaluated these criteria and provided empirical evidence of their efficiency. This study

081 highlights the efficacy of our approach and its 131
 082 potential to enhance text generation by em- 132
 083 ploying diffusion models. In future works, the 133
 084 methodology used in this paper can be devel- 134
 085 oped even further in order to understand better 135
 086 and evaluate newly trained DLMs. 136

087 **2 Related Work**

088 **2.1 Diffusion Language Models**

089 When applied to discrete data, diffusion models 138
 090 have demonstrated promising results in image 139
 091 generation and captioning (Chen et al., 2022). 140

092 Within NLP, diffusion models have also been 141
 093 successfully integrated into sequence-to-sequence 142
 094 tasks (Savinov et al., 2021; Reid et al., 2023; Gong 143
 095 et al., 2023; Yuan et al., 2022; Lin et al., 2022). 144
 096 Despite their performance being on par with non- 145
 097 autoregressive models, most diffusion models em-
 098 ploy an Encoder-Decoder architecture. For uncon-
 099 ditional language modeling, this approach is not
 100 ideal. Although it is possible to modify probabilis-
 101 tic models to accommodate unconditional text gen-
 102 eration, the supporting evidence for their efficacy
 103 in this area is sparse. It is generally acknowledged
 104 that non-autoregressive models are more adept at
 105 dealing with conditional text generation tasks, such
 106 as machine translation, compared to unconditional
 107 text modeling (Gu et al., 2018).

108 The "Diffusion LM" proposed by Li et al. (2022) 146
 109 aimed to establish a generalized LM capable of 147
 110 unconditional sampling. This model was evaluated 148
 111 based on its capability for controlled, classifier- 149
 112 guided text generation. However, it is worth noting 150
 113 that, unlike other pre-trained models, the "Diffu- 151
 114 sion LM" was not trained on large datasets. More- 152
 115 over, its authors did not share any pre-trained 153
 116 weights, making it necessary to train the model 154
 117 from scratch to compare its performance with other 155
 118 methods. Conversely, a significant advantage of the 156
 119 Simplex-based Diffusion Language Model (SSD) 157
 120 and Plaid models (Han et al., 2023; Gulrajani and 158
 121 Hashimoto, 2023) is that they are available in 159
 122 open access and are pre-trained on extensive text 160
 123 datasets. 161

124 Both Self-conditioned Embedding Diffusion 162
 125 (SED) and Continuous Diffusion for Categorical 163
 126 Data (CDCD) have utilized large datasets for pre- 164
 127 training their Diffusion LMs (Strudel et al., 2023; 165
 128 Dieleman et al., 2022). However, neither of them 166
 129 has provided trained model weights or source code. 167
 130

These diffusion models are appealing to use for

comparison due to the different approaches used for 131
 training them. For instance, while CDCD utilizes 132
 a score interpolation objective, SSD works with a 133
 simplex-based method. On the other hand, Plaid is 134
 defined with a Variational Lower Bound objective 135
 (Kingma and Welling, 2014; Kingma et al., 2021). 136

137 **2.2 Early Exiting Methods**

138 The early exit technique is an approach for reducing 139
 139 computational load (Graves, 2016). It especially 140
 140 benefits transformer-based architectures, where in- 141
 141 termediate hidden states maintain consistent shapes 142
 142 across layers. As a result, early exiting has become 143
 143 a standard technique for downstream tasks with pre- 144
 144 trained LMs (Zhou et al., 2020; Liu et al., 2020; 145
 145 Balagansky and Gavrilov, 2022; Gao et al., 2023).

146 **3 Preliminaries**

147 This section will briefly describe the parts essential 148
 148 for understanding various DLMs – CDCD, Plaid, 149
 149 and SSD. While each framework contains many 150
 150 nuances necessary to train and generate samples, 151
 151 we will cover details on loss evaluation and how 152
 152 architecture is defined to evaluate this loss. More 153
 153 concretely, all these models model a categorical dis- 154
 154 tribution over tokens, making it possible to evaluate 155
 155 early exiting with them.

156 For each model, we start with a sequence of 157
 157 tokens $x \in V^l$, where V is a vocabulary set and l 158
 158 is the length of a sequence.

159 These tokens are embedded with the embed- 160
 160 ding matrix $E \in \mathbb{R}^{|V| \times d}$, where d is an embed- 161
 161 ding size, and produce $X_0 \in \mathbb{R}^{l \times d}$. Subscript 0 162
 162 here states that these embeddings do not contain 163
 163 noise, commonly used in the diffusion probabilistic 164
 164 model. Models then operate with the noised em- 165
 165 beddings $X(t)$ where the noise amount depends 166
 166 on timestep t .

167 Also, on top of these models, it is expected to 168
 168 see a layer producing a categorical distribution over 169
 169 possible tokens $p(x|\dots)$, conditioned on arbitrary 170
 170 entities (usually X and t). It is important to note 171
 171 that talking about a distribution over tokens is inter- 172
 172 changeable with discussing a distribution over their 173
 173 corresponding embeddings since each token maps 174
 174 to a specific embedding vector e in \mathbb{R}^d from the em- 175
 175 bedding matrix E , and the process can be reversed. 176
 176 Throughout our discussion, this distribution will 177
 177 also be referred to as $p(e|\dots)$.

3.1 Continuous Diffusion for Categorical Data

Continuous Diffusion for Categorical Data (CDCD) operates with noised input embeddings $\mathbf{X}(t) \in \mathbb{R}^{l \times d}$, where the amount of noise depends on timestep $t \in [0; 1]$. This process predicts a denoised sequence with categorical distribution $p(\mathbf{x}|\mathbf{X}(t), t)$. This distribution is obtained by predicting logits of shape $\mathbb{R}^{l \times |V|}$ and applying the softmax function. Subsequently, cross-entropy loss is applied to estimate $p(\mathbf{x}|\mathbf{X}(t), t)$; i.e.,

$$\mathcal{L}_{CDCD} = -\log(p(\mathbf{x}|\mathbf{X}(t), t)).$$

An estimation of the score function is calculated as $\hat{\mathbf{S}}(\mathbf{X}(t), t) = \frac{\hat{\mathbf{X}}_0(\mathbf{X}(t), t) - \mathbf{X}(t)}{t^2}$, where $\hat{\mathbf{X}}_0(\mathbf{X}(t), t) = \mathbb{E}_{p(\mathbf{x}|\mathbf{X}(t), t)}[[\mathbf{E}, \dots, \mathbf{E}]] := \mathbb{E}_{p(\mathbf{e}|\mathbf{X}(t), t)}[[\mathbf{E}, \dots, \mathbf{E}]] \in \mathbb{R}^{l \times d}$ represents the estimation of the denoised embeddings based on their probabilities (Karras et al., 2022). $\hat{\mathbf{S}}(\mathbf{X}(t), t)$ then could be passed to an arbitrary ODE solver to obtain samples from the model.

3.2 Plaid

Plaid uses simple loss derived from Variational Lower Bound objective (Kingma and Welling, 2014; Kingma et al., 2021; Gulrajani and Hashimoto, 2023)

$$\mathcal{L}_{VLB} = -\frac{1}{2} \mathbb{E}_{t, \mathbf{Z}_t} \left[\frac{d}{dt} \frac{1}{\sigma^2(t)} \|\mathbf{X}_0 - \hat{\mathbf{X}}_0(\mathbf{Z}_t)\|_2^2 \right].$$

Here $\hat{\mathbf{X}}_0(\mathbf{Z}_t) \in \mathbb{R}^{l \times d}$ is an estimation of denoised embeddings from noise \mathbf{Z}_t at time step t , $t \sim U[0; 1]$, $\mathbf{Z}_t \sim q(\mathbf{Z}_t|\mathbf{x})$ is a distribution of forward process defined as $q(\mathbf{Z}_0|\mathbf{x}) = \mathcal{N}(\mathbf{X}_0; \sigma^2(0))$, $q(\mathbf{Z}_t|\mathbf{Z}_s) = \mathcal{N}(\mathbf{Z}_s; \sigma^2(t) - \sigma^2(s))$. $\sigma^2(t)$ is modelled following Kingma et al. (2021).

Notably, while $\hat{\mathbf{X}}_0$ could be modeled in continuous space, doing so will require a model to remember initial embeddings, which is redundant. Instead, plaid uses a categorical reparametrization similar to one used with Section 3.1 and directly learns $p(\mathbf{e}|\mathbf{Z}_t)$ to estimate $\hat{\mathbf{X}}_0(\mathbf{Z}_t) = \mathbb{E}_{p(\mathbf{e}|\mathbf{Z}_t)}[[\mathbf{E}, \dots, \mathbf{E}]]$.

3.3 Simplex-based Diffusion Language Model

Simplex-based Diffusion Language Model (SSD) is a third Diffusion LM tested with unconditional text generation for general language modeling.

Starting with token sequence \mathbf{x} , SSD firstly defines the operation for almost-one-hot encodings

of \mathbf{x} , namely *logits generation*. For token \mathbf{x}_i , its continuous representation is defined as $\mathbf{X}_{i,j} = K$ if $\mathbf{x}_i = \mathbf{V}_j$, and $\mathbf{X}_{i,j} = -K$ otherwise. $K \in \mathbb{R}^+$ here is a hyperparameter, and \mathbf{V} is a vocabulary.

Then, for the forward process of the diffusion with sequence \mathbf{x} , we evaluate logit generation for tokens $\mathbf{x}_{c:l}$, where c is a context length. Noise is progressively added to almost-one-hot encodings of the text $\mathbf{X}_{c:l}(t)$, leading to normal distribution across logits at the end.

The model is then trained with a loss

$$\mathcal{L}_{SSD} = \mathbb{E}_{c,t} \left[\sum_{j=c}^{l-1} -\log(p(\mathbf{x}_j|\mathbf{X}_{c:l}(t), \mathbf{x}_{<c})) \right],$$

where $c \sim U(1; l)$ is the length of context for a generation, and $t \sim U(1; T)$, and $p(\mathbf{x}_j|\mathbf{X}_{c:l}(t), \mathbf{x}_{<c})$ is a categorical distribution over tokens from vocabulary.

4 Early Exiting with DLMS

While CDCD, Plaid, and SSD define different views on training DLMS, they share a similarity in how distribution on denoised text is defined. More concretely, they all define a categorical distribution over possible embeddings (and thus over possible tokens).

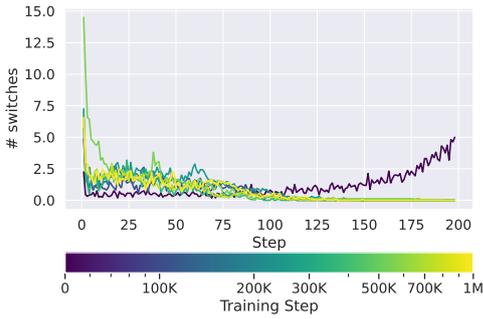
This fact leads to the question of *how the distribution of possible tokens changes with time*.

4.1 Emergence of Early Exiting Behavior

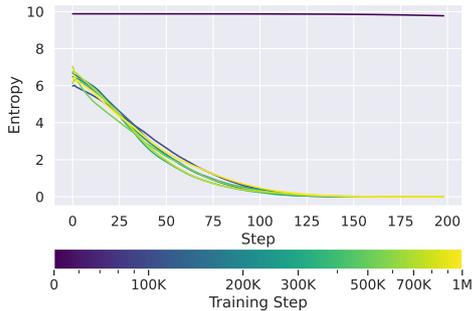
To explore token behavior during generation, we analyze the number of token switches (changes in tokens after each generation step) in CDCD. Note that Dieleman et al. (2022) did not release pre-trained models or training codes for CDCD. To perform experiments with this framework, we reproduced it and trained our model, namely **Democratized Diffusion Language Model** (DDLMS). See Appendix Section A for reproduction details.

We evaluate token switches at different pre-training checkpoints and at each time step t during generation for DDLMS. Additionally, we examine the entropy of the embedding prediction $p(\mathbf{x}|\mathbf{X}(t), t)$. Sequences with 200 steps are sampled for this analysis (see Figure 1).

Interestingly, the model shows zero token switches after approximately the 100th sampling step. This suggests a potential for adaptive early exiting in DDLMS generation since, for nearly half of



(a)



(b)

Figure 1: (a) The number of token switches and (b) the entropy of $p(\mathbf{x}|\mathbf{X}(t), t)$. Color represents the training step, while the x-axis is the diffusion generation step. The trained model reaches the minimum entropy value before the generation process ends, and the resulting samples remain unchanged. See Section 4.1 for more details.

the generation steps, the sampling algorithm only made minor adjustments to predicted embeddings without changing the generated tokens. Depending on the sequence, adaptive early exiting will make it possible to dynamically evaluate when we can halt the generation process, potentially greatly reducing the computations needed for sampling.

To understand why the trained model tends towards minimal token switches early on in the generation process, we examined the L2 norm of $\hat{\mathbf{X}}_0$ and \mathbf{X} during generation¹ (refer to Figure 2). We found that $\hat{\mathbf{X}}_0$ rapidly reaches an L2 norm of 16, the L2 norm of normalized embeddings during

¹For the reader’s convenience, it is essential to remember that \mathbf{X} are embeddings passed to the model as an input. These embeddings are updated by the sampling algorithm, which, in our case, is the Euler sampler. At the same time, $\hat{\mathbf{X}}_0$ are embeddings produced by the model to estimate the score function. These embeddings and their statistics differ during the generation process: $\hat{\mathbf{X}}_0$ could change fast, while \mathbf{X} will change slowly.

Noise	AR-NLL	dist ₁	dist ₂	dist ₃	s.-BLEU
0.0	0.44	0.00	0.00	0.00	1.00
0.5	3.10	0.24	0.47	0.60	0.58
0.8	3.50	0.41	0.74	0.84	0.47
0.9	3.62	0.48	0.83	0.92	0.49
1.0	3.72	0.49	0.86	0.94	0.48
1.1	3.86	0.51	0.88	0.90	0.47
1.2	4.01	0.52	0.89	0.95	0.44

Table 1: Performance of DDLM depending on the initial noise scale of \mathbf{X} . Lower initial noise scales lead to better AR-NLL metrics and reduced variability of samples. See Section 4 for more details.

pre-training. This aligns with our observation of the entropy of $p(\mathbf{x}|\mathbf{X}, t)$ reaching near-zero values within 100 generation steps. Fascinatingly, the L2 norm of \mathbf{X} first reduces and then increases from its large initialization value, suggesting that \mathbf{X} travels from one point on the embedding sphere surface to another via its interior.

To support this hypothesis, we evaluate the cos between score $\hat{\mathbf{S}}$ with final score $\hat{\mathbf{S}}_0$, and the cos between \mathbf{X} with final \mathbf{X}_0 during the generation process. After the 100th step, the scoring angle stops changing, indicating that the model settles on the final embedding improvement direction of mid-generation. This constant direction forces \mathbf{X} to the embedding sphere boundary, leading to high-confidence results and near-zero token switches.

Empirical evidence suggests that \mathbf{X} traverses between two points on the surface of a sphere via its interior. By reducing the initial noise scale, we can adjust the trajectory of \mathbf{X} . See Figure 3 and Table 1 for our results. We find that a lower initial noise scale makes it possible for $\|\mathbf{X}\|_2$ to reach its minimum value during generation more quickly. However, this approach limits the variability of samples. While our findings show that using a noise scale of 0.9 is optimal, we will use a scale of 1.0 in later experiments for convenience.

4.2 Exploring Early Exit Criteria

The concept of early exiting is a well-established practice in various research fields of Deep Learning (Graves, 2016; Liu et al., 2020; Zhou et al., 2020; Balagansky and Gavrilo, 2022; Graves, 2016). Consequently, there are numerous methods available for performing an early exit.

Entropy criterion, described by Liu et al. (2020), is one of the most common early exit techniques. This method performs an exit when entropy drops below a certain threshold. A major down-

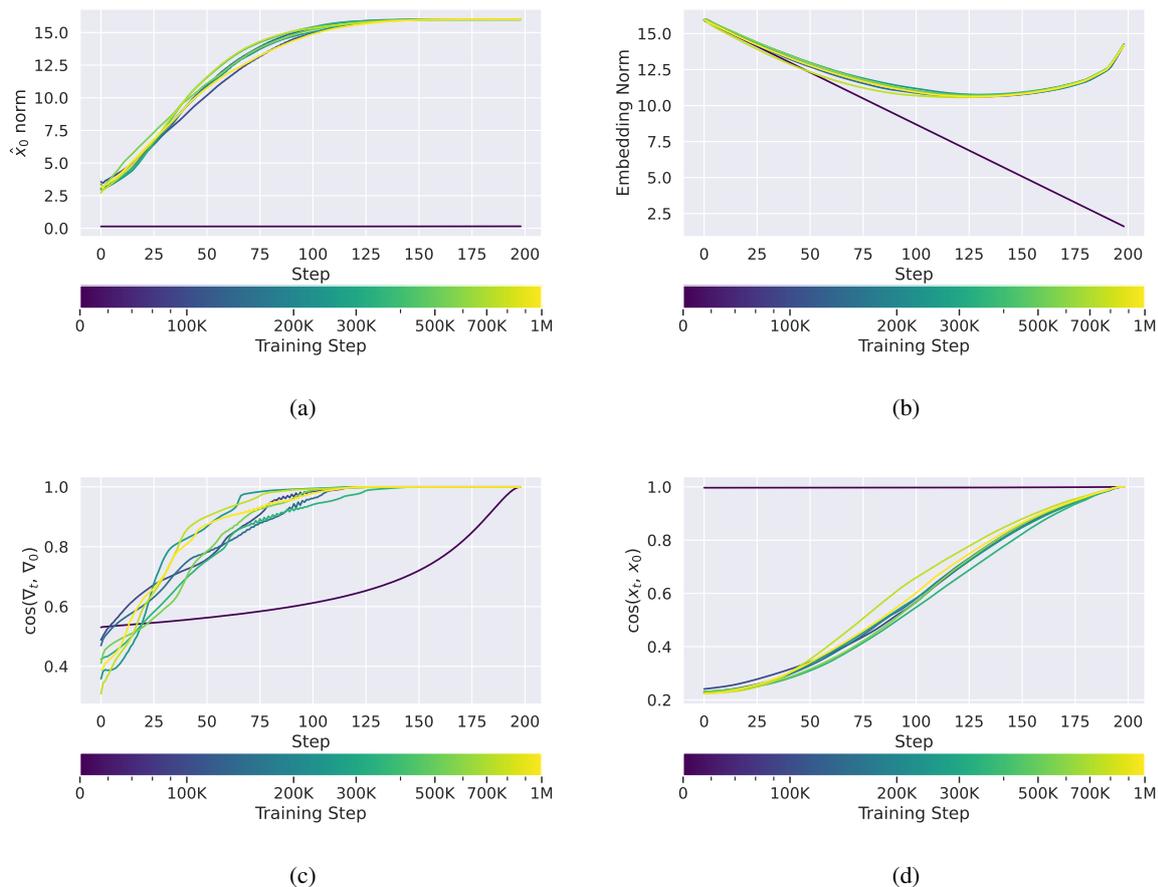


Figure 2: (a) The L2 norm of embeddings $\|\hat{\mathbf{X}}_0\|_2$, (b) the L2 norm of embeddings $\|\mathbf{X}\|_2$, (c) \cos of the angle between score estimation $\hat{\mathbf{S}}$ and final score in the end of generation, and (d) \cos of the angle between embedding x and final embedding in the end of generation. Color represents the training step, while the x-axis is the diffusion generation step. See Section 4.1 for more details.

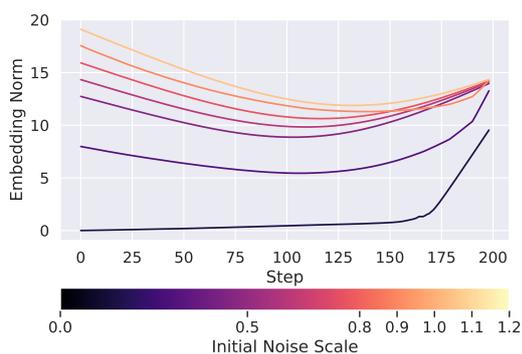


Figure 3: The L2 norm of embeddings $\|\mathbf{X}\|_2$ during the generation process for different initial scales of $\|\mathbf{X}\|_2$ for DDLM. Color represents the initial noise scale, while the x-axis is the diffusion generation step. See Section 4 for more details.

side of the entropy criterion is that it disregards the output dynamics, resulting in overly confident classifiers. Refer to Algorithm 1 for more details.

Patience-based criterion, as proposed by (Zhou et al., 2020), addresses the limitations of the Entropy criterion. It is formulated as follows: if the classifier predictions remain unchanged for a series of t consecutive steps, the model initiates an exit. A notable drawback of Patience is its insensitivity to the scale of the changes. It can trigger an exit due to minor alterations in the output distribution or persist even when significant changes occur. Another drawback of this approach is that it requires a substantial number of steps for the patience value to become meaningful, which is not ideal when the goal is to minimize the number of steps. An exit criterion based on the count of token switches during generation can be seen as Patience-based, as it terminates generation when the number of al-

319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337

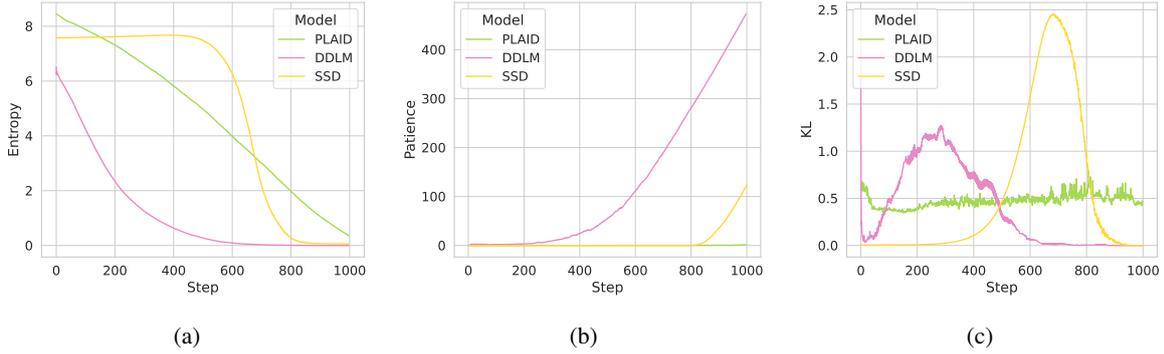


Figure 4: (a) Entropy, (b) unchanged step count, and (c) KL-Divergence are used for different criteria in DDLM, SSD, and Plaid. Generation is halted when the threshold values are met. DDLM reaches the threshold early on, while SSD does so in later stages. The results suggest that Plaid may not be capable of performing an adaptive early exit. See Section 4.2 for more details.

338 tered tokens falls below the threshold value for a
 339 sequence of generation steps. Further details are
 340 provided in Algorithm 2.

341 **KL criterion** overcomes the drawbacks of the
 342 Patience-based criterion (Gao et al., 2023). This
 343 criterion triggers an exit when the KL Divergence
 344 between the current diffusion step’s distribution
 345 and the previous one falls below a certain thresh-
 346 old. This approach reduces the required number
 347 of steps by half and enhances the quality of the
 348 generated texts, which we demonstrate later. Refer
 349 to Algorithm 3 for more details.

350 As seen in Figure 4, all the criteria applied to
 351 DDLM show that it may be possible to halt sam-
 352 pling during generation. For SSD, these criteria
 353 suggest stopping after the 800th step out of 1000.
 354 On the other hand, for Plaid, we observed that
 355 entropy decayed linearly during generation while
 356 other criteria remained constant. This suggests the
 357 possibility of Plaid performing poorly with adap-
 358 tive early exiting methods.

359 We aim to see how these early exiting strategies
 360 perform when applied to various DLMs.

361 4.3 Optimal Number of Steps

362 In this experiment, we want to compare different
 363 adaptive early exiting criteria to the fixed early
 364 exiting strategy on three baseline models: DDLM,
 365 Plaid, and SSD. For each model, we aim to find
 366 the criteria and the corresponding thresholds that
 367 both reduce the mean amount of observed steps
 368 and produce high-quality samples.

369 To evaluate sample quality, we analyze several
 370 adaptive early exiting criteria compared to a fixed

371 early exiting strategy at specific steps. We evaluate
 372 all models in the Prefix-32 setup with 1000 gen-
 373 eration steps. For each generation step, we assess
 374 the AR-NLL metric. Based on results from Sec-
 375 tion 4.2, we expect DDLM to perform an early exit
 376 around the 600th generation step. For SSD, we ex-
 377 pect to see adaptive early exiting capabilities after
 378 the 800th step. Meanwhile, we do not expect adap-
 379 tive early exiting for Plaid since entropy reaches its
 380 minimum only at the end of the generation process.
 381 However, it may be possible for Plaid to perform
 382 early exiting with a fixed exit step.

383 See Figure 5 for results. As hypothesized, we
 384 observed that DDLM could perform adaptive early
 385 exiting during generation after the 600th step. Fur-
 386 thermore, the KL criterion allowed us to perform an
 387 earlier exit than the fixed criterion for a fixed AR-
 388 NLL value. More specifically, we exited 50 steps
 389 earlier on average without losing sample quality.

390 For the SSD model, early exiting with the KL
 391 criterion also performed marginally better than the
 392 fixed strategy. The computation gain for this model
 393 was around 10 steps. KL, Patience, and Fixed cri-
 394 teria showed comparable performance and allowed
 395 an early exit at the 850th step without any loss
 396 in quality compared to the final sample from the
 397 1000th step.

398 As we initially hypothesized, we did not observe
 399 adaptive early exiting capabilities in Plaid. Both
 400 Patience and KL criteria largely underperformed
 401 when compared to fixed and Entropy criteria. At
 402 the same time, the Entropy criterion does not dis-
 403 play an advantage over the fixed exit criterion. This
 404 result aligns with our observation of the values of

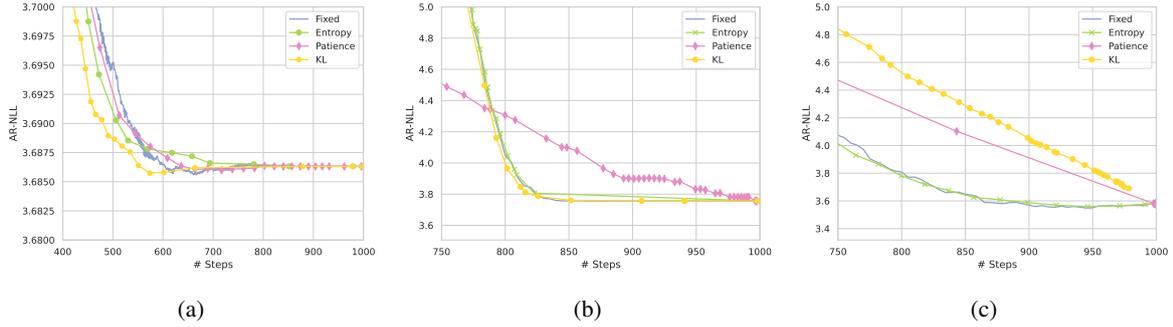


Figure 5: (a) AR-NLL for the different exit criteria with DDLML, (b) SSD, and (c) Plaid with 1k samples of the C4 validation set. See Section 4.3 for more details.

various criteria during generation, where only the Entropy criterion provided meaningful information regarding the sampling process dynamic. However, even though adaptive early exiting did not succeed with Plaid, we observed that AR-NLL stopped changing with fixed criterion after the 900th generation step. This suggests that early exiting can still be performed to reduce computational footprint during generation.

Our results show a speed increase of 40% for DDLML, 20% for SSD, and 10% for Plaid. This enables us either to generate text faster or improve text quality by allowing more steps in the same time frame. We also observed that early exiting methods do not hurt the diversity of samples² (see Figure 6).

See Appendix Figure 8 for results with samples of length 256.

4.4 On Convergence of Early Exiting Methods

We evaluate the sample dynamics during generation with GPT-4 (OpenAI, 2023) to understand the sample dynamics during generation. Recently, Rafailov et al. (2023) showed that this approach is comparable to human judgment and helps assess many samples for different time steps. We also calculate the Word Error Rate (WER) score between samples during generation and the sample from the final step.

With such side-by-side assessment, our end goal is to understand the convergence of generations. GPT-4 allows us to compare samples with reference

²One may find this result to contradict one observed with Section 4 and Table 1. However, for experiments with noise scales, reduced variability is observed for small initial noise scales, leading to deterministic generation. At the same time, a noise scale equal to 1.0 produces diverse samples, while early exiting methods do not hurt this variability.

texts by considering their semantics, thus providing a broader evaluation. Meanwhile, WER shows the differences at the word level. See Appendix Section B for more details on GPT-Score.

Our results are presented in Figure 7. DDLML converged with GPT-Score after the 600th step, and there was no variance in samples afterward. For SSD, we observed the same behavior after the 850th step. Meanwhile, for Plaid, we did not observe any convergence after the 900th step with GPT-Score, and the GPT-Score of the side-by-side comparison with the final sample was large enough. The GPT-4 response indicated minor differences with the reference text, while WER reached low values, indicating that a fixed early exit could still be performed despite entropy not reaching its minimum. See Appendix Section C for sample examples.

5 Discussion

Early Exiting Strategies. One notable observation is that for both CDCD and SSD models, we can effectively implement adaptive techniques that allow the generation process to stop prematurely. In contrast, the Plaid model can halt generation without such adaptiveness. Most importantly, employing these early exiting tactics does not result in a decline in the generated content quality.

This finding has dual benefits. It can quicken text generation without quality loss or increase generation steps within a fixed timeframe to improve output quality. These enhancements promise broader adoption and ongoing advancement of DLMs.

Identifying Issues in DLMs. The ability to stop the text generation process early also signals opportunities to refine DLM design. We contemplate two possibilities: a) varying computational needs for

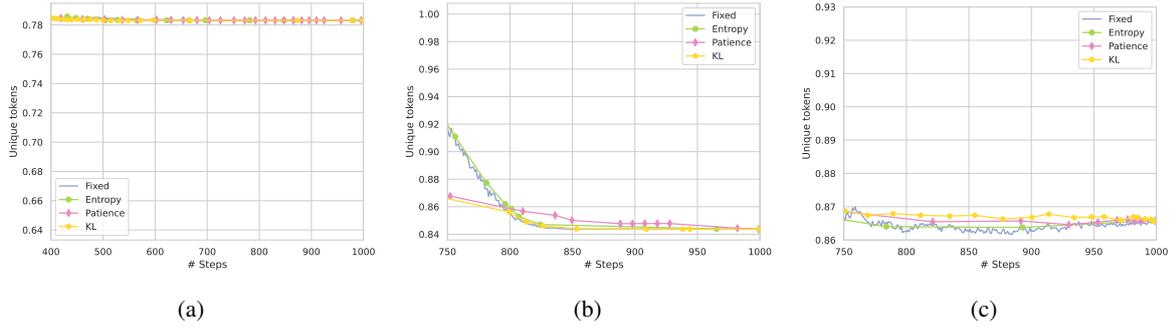


Figure 6: Fraction of unique tokens for the different exit criteria with (a) DDLM, (b) SSD, and (c) Plaid with 1k samples of the C4 validation set. Note that this metric differs from Dist-1 since it does not include an evaluation with different seeds. See Section 4.3 for more details.

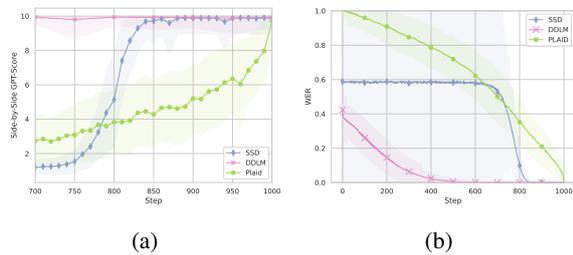


Figure 7: (a) Side-by-side GPT-Score and (b) WER with final sample for DDLM, SSD, and Plaid models with a fixed early exiting mechanism. The plot is truncated to 400 generation steps for GPT-Score. DDLM stabilizes at step 600, SSD at 850, and Plaid continues evolving until the end. However, after step 900, Plaid shows minimal WER differences. For further information, refer to Section 4.3.

different text generation tasks suggest early halting is apt for simple texts to prevent over-processing and beneficial for complex texts for additional computation; b) the computational effort may not vary with text complexity, suggesting that the capacity for early halting could point to design inefficiencies in DLMs (i.e., early exiting should not occur for properly trained and designed DLMs, thus indicating on issues with existing models). In the latter case, if the emergence of an early exit is an issue in the design of current DLMs, our research is a valuable methodology tool to evaluate and probe the performance of new pre-trained models.

Considering dynamic generation processes is vital for deeply understanding model capabilities and their constraints. Such dynamic evaluations are often overlooked, with many studies preferring to assess a model’s static performance using metrics like data likelihood (Gulrajani and Hashimoto,

2023). However, lessons from the Computer Vision field show that examining process dynamics can yield rich insights into specific cases (Karras et al., 2023).

Directions for Future Research. Our methodology offers insights into assessing the performance of emerging DLMs, noting that the option for early exiting could indicate underlying issues in the trained models. Therefore, future investigations could build upon our approach, incorporating new evaluation criteria or exploring DLMs that do not support early exiting. This could shed more detail on the strengths and potential weaknesses of these models.

6 Limitations

This paper only used our re-implementation of DLM trained with the CDCD framework, SSD, and Plaid models. We omitted other Diffusion Models, such as GENIE or DiffuSeq (Lin et al., 2022; Gong et al., 2023), since there is no evidence that these frameworks can perform unconditional text generation if trained in such a manner.

Our experiments involve our own DDLM model, which a reproduction of DLM trained with the CDCD framework. It is not a precise reproduction, as there is no source code available for CDCD. Nevertheless, we believe that conducting experiments on our model, which was trained with a score interpolation objective, made it possible for us to present more comprehensive results in this paper.

Our analysis focuses on the AR-NLL metric to frequently evaluate models during generation. However, our evaluation with GPT-4 indicates that no issues with the analysis should have occurred,

526	and our baseline models converged during generation.	581
527		582
528	References	
529	Nikita Balagansky and Daniil Gavrilov. 2022. Palbert: Teaching albert to ponder . In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 14002–14012. Curran Associates, Inc.	583
530		584
531		585
532		586
533	Sid Black, Gao Leo, Phil Wang, Connor Leahy, and Stella Biderman. 2021. GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow . If you use this software, please cite it using these metadata.	587
534		
535		
536		
537		
538	Ting Chen, Ruixiang Zhang, and Geoffrey Hinton. 2022. Analog bits: Generating discrete data using diffusion models with self-conditioning .	588
539		589
540		
541	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek B Rao, Parker Barnes, Yi Tay, Noam M. Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Benton C. Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier García, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Díaz, Orhan Firat, Michele Catasta, Jason Wei, Kathleen S. Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. Palm: Scaling language modeling with pathways . <i>ArXiv</i> , abs/2204.02311.	590
542		591
543		592
544		593
545		594
546		595
547		596
548		597
549		598
550		599
551		600
552		601
553		602
554		603
555		604
556		605
557		606
558		607
559		608
560		609
561		610
562		611
563		612
564	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.	613
565		614
566		615
567		616
568		617
569		618
570		619
571		620
572		621
573	Sander Dieleman, Laurent Sartran, Arman Roshanai, Nikolay Savinov, Yaroslav Ganin, Pierre H. Richemond, Arnaud Doucet, Robin Strudel, Chris Dyer, Conor Durkan, Curtis Hawthorne, Rémi Leblond, Will Grathwohl, and Jonas Adler. 2022. Continuous diffusion for categorical data .	622
574		623
575		624
576		625
577		626
578		627
579	Xiangxiang Gao, Wei Zhu, Jiasheng Gao, and Congrui Yin. 2023. F-pabee: Flexible-patience-based early exiting for single-label and multi-label text classification tasks .	628
580		629
		630
		631
		632
		633
		634
		635
		636
		637
		638
		639
		640
		641
		642
		643
		644
		645
		646
		647
		648
		649
		650
		651
		652
		653
		654
		655
		656
		657
		658
		659
		660
		661
		662
		663
		664
		665
		666
		667
		668
		669
		670
		671
		672
		673
		674
		675
		676
		677
		678
		679
		680
		681
		682
		683
		684
		685
		686
		687
		688
		689
		690
		691
		692
		693
		694
		695
		696
		697
		698
		699
		700
		701
		702
		703
		704
		705
		706
		707
		708
		709
		710
		711
		712
		713
		714
		715
		716
		717
		718
		719
		720
		721
		722
		723
		724
		725
		726
		727
		728
		729
		730
		731
		732
		733
		734
		735
		736
		737
		738
		739
		740
		741
		742
		743
		744
		745
		746
		747
		748
		749
		750
		751
		752
		753
		754
		755
		756
		757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800

635	of language representations. In <i>International Conference on Learning Representations</i> .	Machel Reid, Vincent Josua Hellendoorn, and Graham Neubig. 2023. <i>DiffusER: Diffusion via edit-based reconstruction</i> . In <i>The Eleventh International Conference on Learning Representations</i> .	688
636			689
637	Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori B. Hashimoto. 2022. <i>Diffusion-lm improves controllable text generation</i> .		690
638			691
639			
640	Zhenghao Lin, Yeyun Gong, Yelong Shen, Tong Wu, Zhihao Fan, Chen Lin, Weizhu Chen, and Nan Duan. 2022. <i>Genie: Large scale pre-training for text generation with diffusion model</i> .	Nikolay Savinov, Junyoung Chung, Mikolaj Binkowski, Erich Elsen, and Aaron van den Oord. 2021. <i>Step-unrolled denoising autoencoders for text generation</i> .	692
641			693
642			694
643			
644	Weijie Liu, Peng Zhou, Zhiruo Wang, Zhe Zhao, Haotang Deng, and Qi Ju. 2020. <i>FastBERT: a self-distilling BERT with adaptive inference time</i> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 6035–6044, Online. Association for Computational Linguistics.	Jiaming Song, Chenlin Meng, and Stefano Ermon. 2020. <i>Denoising diffusion implicit models</i> . <i>arXiv:2010.02502</i> .	695
645			696
646			697
647			
648			
649			
650			
651	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. <i>Roberta: A robustly optimized bert pretraining approach</i> . Cite arxiv:1907.11692.	Robin Strudel, Corentin Tallec, Florent Altché, Yilun Du, Yaroslav Ganin, Arthur Mensch, Will Sussman Grathwohl, Nikolay Savinov, Sander Dieleman, Laurent Sifre, and Rémi Leblond. 2023. <i>Self-conditioned embedding diffusion for text generation</i> .	698
652			699
653			700
654			701
655			702
656	Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. 2022. <i>GLIDE: towards photorealistic image generation and editing with text-guided diffusion models</i> . In <i>International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA</i> , volume 162 of <i>Proceedings of Machine Learning Research</i> , pages 16784–16804. PMLR.	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. <i>Attention is all you need</i> . In <i>Advances in Neural Information Processing Systems</i> , pages 5998–6008.	703
657			704
658			705
659			706
660			707
661			
662			
663			
664			
665	OpenAI. 2023. <i>Gpt-4 technical report</i> . <i>ArXiv</i> , abs/2303.08774.	Hongyi Yuan, Zheng Yuan, Chuanqi Tan, Fei Huang, and Songfang Huang. 2022. <i>Seqdiffuseq: Text diffusion with encoder-decoder transformers</i> .	708
666			709
667			710
668			
669			
670	Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron C. Courville. 2018. <i>Film: Visual reasoning with a general conditioning layer</i> . In <i>AAAI</i> .	Wangchunshu Zhou, Canwen Xu, Tao Ge, Julian McAuley, Ke Xu, and Furu Wei. 2020. <i>Bert loses patience: Fast and robust inference with early exit</i> . In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 18330–18341. Curran Associates, Inc.	711
671			712
672			713
673			714
674			715
675			716
676			
677			
678	Krishna Pillutla, Swabha Swayamdipta, Rowan Zellers, John Thickstun, Sean Welleck, Yejin Choi, and Zaid Harchaoui. 2021. <i>Mauve: Measuring the gap between neural text and human text using divergence frontiers</i> . In <i>NeurIPS</i> .		
679			
680			
681			
682	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. <i>Language models are unsupervised multitask learners</i> .		
683			
684			
685			
686			
687			
688			
689			
690			
691			
692			
693			
694			
695			
696			
697			
698			
699			
700			
701			
702			
703			
704			
705			
706			
707			
708			
709			
710			
711			
712			
713			
714			
715			
716			

Algorithm 1: Entropy algorithm

Require: Diffusion model $f_\theta(\cdot, \cdot)$, entropy threshold e_t , maximum number of diffusion steps N_{\max} , timestamps array t .

- 1: step $\leftarrow 0$
- 2: $x \leftarrow X \sim \mathcal{N}(0, I)$
- 3: **while** step $< N_{\max}$ **do**
- 4: $p(\text{tokens}_{\text{cur}}), \hat{x} \leftarrow f_\theta(x, t[\text{step}])$
- 5: $e \leftarrow \text{entropy}(p(\text{tokens}_{\text{cur}}))$
- 6: **if** $e \geq e_t$ **then**
- 7: **return** $p(\text{tokens}_{\text{cur}})$
- 8: **end if**
- 9: $x \leftarrow \text{Euler}(x, \hat{x}, t)$
- 10: step $\leftarrow \text{step} + 1$
- 11: **end while**
- 12: **return** $p(\text{tokens}_{\text{cur}})$

717

Algorithm 2: Patience algorithm

Require: Diffusion model $f_\theta(\cdot, \cdot)$, patience threshold p , maximum number of diffusion steps N_{\max} , timestamps array t

- 1: $\text{step} \leftarrow 0$
- 2: $p_{\text{cur}} \leftarrow 0$
- 3: $x \leftarrow X \sim \mathcal{N}(0, I)$
- 4: **while** $\text{step} < N_{\max}$ **do**
- 5: $p(\text{tokens}_{\text{cur}}), \hat{x} \leftarrow f_\theta(x, t[\text{step}])$
- 6: $\text{tokens}_{\text{cur}} \leftarrow \text{argmax}(p(\text{tokens}_{\text{cur}}))$
- 7: **if** $\text{step} > 0$ **then**
- 8: **if** $\text{tokens}_{\text{cur}} = \text{tokens}_{\text{prev}}$ **then**
- 9: $p_{\text{cur}} \leftarrow p_{\text{cur}} + 1$
- 10: **else**
- 11: $p_{\text{cur}} \leftarrow 0$
- 12: **end if**
- 13: **if** $p_{\text{cur}} \geq p$ **then**
- 14: **return** $p(\text{tokens}_{\text{cur}})$
- 15: **end if**
- 16: **end if**
- 17: $x \leftarrow \text{Euler}(x, \hat{x}, t)$
- 18: $\text{tokens}_{\text{prev}} \leftarrow \text{tokens}_{\text{cur}}$
- 19: $\text{step} \leftarrow \text{step} + 1$
- 20: **end while**
- 21: **return** $p(\text{tokens}_{\text{cur}})$

Algorithm 3: KL algorithm

Require: Diffusion model $f_\theta(\cdot, \cdot)$, divergence-threshold d , maximum number of diffusion steps N_{\max} , parameter $\text{min_steps} \approx 0.25N_{\max}$, timestamps array t ,

- 1: $\text{step} \leftarrow 0$
- 2: $x \leftarrow X \sim \mathcal{N}(0, I)$
- 3: **while** $\text{step} < N_{\max}$ **do**
- 4: $p(\text{tokens}_{\text{cur}}), \hat{x} \leftarrow f_\theta(x, t[\text{step}])$
- 5: **if** $\mathcal{D}(p(\text{tokens}_{\text{cur}}) || p(\text{tokens}_{\text{prev}})) > d_t$
and $s \geq \text{min_steps}$ **then**
- 6: **return** $p(\text{tokens}_{\text{cur}})$
- 7: **end if**
- 8: $x \leftarrow \text{Euler}(x, \hat{x}, t)$
- 9: $\text{step} \leftarrow \text{step} + 1$
- 10: $p(\text{tokens}_{\text{prev}}) \leftarrow p(\text{tokens}_{\text{cur}})$
- 11: **end while**
- 12: **return** $p(\text{tokens}_{\text{cur}})$

A Reproducing CDCD

Comparing the CDCD model with other diffusion models is an intriguing challenge due to its unique objectives that set it apart from conventional DLMs. However, the lack of a publicly available training code for the CDCD limits such research. Therefore, we have reproduced this model in order to understand the differences between CDCD and other frameworks. We will briefly describe the essential parts of the CDCD framework and then go into detail about our reproduction of the CDCD.

A.1 Understanding CDCD Framework

Once loss and score functions are defined, CDCD implies several details must be considered before training a model.

The first of them is **Embeddings normalization**. As the model with the \mathcal{L}_{CDCD} loss function is forced to distinguish correct embeddings from noisy ones, a naive application of such an objective will lead to uncontrollable growth of the embeddings norm to make them easier to distinguish. CDCD applies L_2 normalization during training to prevent an uncontrolled growth of embedding norms.

Second, the score interpolation objective implies sampling the time t from some distribution during the training. While it is possible to sample t uniformly in $[0; 1]$, Dieleman et al. (2022) used **Time Warping** method. Dieleman et al. (2022) trained CDF of time $F_\phi(t)$ following Kingma et al. (2021). More concretely, for the CDCD framework, $F_\phi(t)$ is trained with a loss $\mathcal{L}_{TW} = \|\tilde{F}_\phi(t) - \mathcal{L}_{CDCD}(\dots, t)\|$, where $\tilde{F}_\phi(t)$ is the unnormalized CDF parametrized with ϕ . We can obtain samples from it by normalizing and inverting $\tilde{F}_\phi(t)$. $p(x|\mathbf{X}, t)$ is then conditioned on t via conditional layer normalization (Perez et al., 2018).

Finally, since our model is trained akin to Masked Language Models to fill noisy tokens with real ones, it is essential to define the mechanism to select specific tokens to inject noise, i.e., **Noise masking**. The first approach, prefix masking, involves injecting noise into the embedding sequence continuation while keeping its beginning intact. Alternatively, noise can be injected at random sequence positions, similar to Masked Language Models training (MLM masking) (Devlin et al., 2019; He et al., 2020; Liu et al., 2019; Lan et al., 2020). The third approach combines the previous

Model	Steps	Sampler	AR-NLL	Dist-1	Dist-2	Dist-3	MAUVE	Zipf's Coef.
Data	N/A	N/A	3.31	N/A	N/A	N/A	N/A	0.90
Prefix-32								
DDLML, 147M	50	Euler	3.72	0.53	0.85	0.90	0.80	0.96
	200		3.65	0.54	0.84	0.90	0.82	0.96
	1000		3.63	0.54	0.84	0.90	0.81	0.96
Plaid, 1.3B	200	DDPM	3.69	0.66	0.88	0.90	0.93	0.86
	500		3.64	0.65	0.87	0.89	0.89	0.87
	1000		3.65	0.65	0.87	0.90	0.94	0.87
SSD, 400M	200	Simplex	4.00	0.66	0.91	0.83	0.82	0.88
	1000		3.75	0.63	0.91	0.83	0.85	0.90
GPT-2, 124M	N/A	N/A	3.21	0.58	0.86	0.89	0.86	0.96
GPT-Neo, 125M	N/A	N/A	3.20	0.60	0.85	0.88	0.83	0.96
Unconditional								
DDLML, 147M	50	Euler	3.98	0.50	0.85	0.93	N/A	1.19
	200		3.77	0.50	0.84	0.92	N/A	1.17
	1000		3.67	0.49	0.83	0.91	N/A	1.16
Plaid, 1.3B	200	DDPM	3.83	0.66	0.92	0.94	N/A	0.93
	500		3.73	0.65	0.91	0.94	N/A	0.94
	1000		3.69	0.65	0.91	0.94	N/A	0.94
SSD, 400M	200	Simplex	6.45	0.57	0.91	0.83	N/A	0.99
	1000		6.55	0.57	0.91	0.83	N/A	1.12
GPT-2, 124M	N/A	N/A	2.62	0.67	0.90	0.90	N/A	1.10
GPT-Neo, 125M	N/A	N/A	2.27	0.66	0.88	0.89	N/A	1.05

Table 2: Evaluation of DDLML, SSD, Plaid, GPT-2, and GPT-Neo with 5k samples of the C4 validation set with the Unconditional and Prefix-32 tasks. The best result across DLMs is bolded. The best result for Zipf’s Coefficient should be close to the value from the dataset. See Section A for more details.

two, injecting noise into random positions in a sequence continuation (mixed masking). The cross-entropy loss \mathcal{L}_{CDCD} is calculated only with noised embeddings.

CDCD is implemented as Transformer (Vaswani et al., 2017). Once all objective embeddings necessary for score interpolation are concatenated, they are passed through Transformer layers to obtain $p(\mathbf{x}|\mathbf{X}, t)$.

A.2 Training DDLML

Following the information provided on the CDCD framework, we trained our version of it, namely the Democratized Diffusion Language Model (DDLML)³. We trained this model using the C4 dataset (Rafel et al., 2020) with 147M parameter models and a sequence length of 64 tokens.

The tokenized training data consisted of a vocabulary $|V| = 32\text{k}$, and the tokens used 256-sized embeddings following Dieleman et al. (2022). We trained DDLML using 8 Nvidia A100 SXM4 80GB

³“Democratized” in the model name stands for the open availability of this model for other researchers.

GPUs, completing one million training steps over approximately 1.5 days. The details on the hyperparameters used can be found in Table 7.

For validation, we extracted 5k examples from the C4 validation set and generated 5 separate continuations using different seeds. Our evaluation of DDLML was carried out in two setups: Unconditional and Prefix-32, where text was generated using a prefixed prompt of 32 tokens in length. We utilized several metrics to assess the quality of the text, including AR-NLL as measured by GPT-Neo-1.3B (Black et al., 2021), the MAUVE metric (Pillutla et al., 2021), average distinct N-grams over 5 samples with a single prompt (where available), and Zipf’s coefficient over token distribution. These metrics cover the various properties of generated texts and make it possible for us to perform an in-depth evaluation of DLMs by evaluating DLMs at each generation step.

While Dieleman et al. (2022) states that small values of t_{\max} can lead to trivial solutions for the score interpolation objective, we hypothesize that applying several normalizations during training,

812 such as normalizing embeddings and noised em- 863
813 beddings, can prevent trivial solutions from emerg- 864
814 ing. 865

815 Additionally, our interest extended to delving 866
816 deeper into noise-masking strategies. While [Diele- 867
817 man et al. \(2022\)](#) favored mixed masking, we sug- 868
818 gested an extension of prefix masking, a component 869
819 of mixed masking, to span masking ([Strudel et al., 870
820 2023](#)). In span masking, a sequence of tokens is di- 871
821 vided into k segments (k being a randomly chosen 872
822 integer between 1 and a fixed constant $k_{max} = 9$) 873
823 by randomly selecting $k - 1$ indices. These in- 874
824 dices define k spans, each subjected to noise with 875
825 a probability of 50%. It is important to note that 876
826 our experimentation with the span masking strategy 877
827 was not aimed at achieving superior performance 878
828 compared to other methods, but rather at uncover- 879
829 ing their distinctions. 880

830 We trained models with different t_{max} values, in- 881
831 cluding $t_{max} \in [10, 50, 300]$. Both models with and 882
832 without time warping were trained for each t_{max} 883
833 value. Furthermore, all these experiments were 884
834 conducted using three masking strategies: MLM, 885
835 prefix, and span. 886

836 For the detailed results of Unconditional, Prefix- 887
837 32, and Enclosed-32 generation, refer to Table 3 888
838 and Appendix Tables 4, 5, and 6. We observed that 889
839 training models with high t_{max} values led to poor 890
840 results with repetitive samples. Comprehensive 891
841 samples were only achieved when t_{max} was reduced 892
842 to 10. Notably, while larger t_{max} values resulted in 893
843 poor samples, the loss values for such setups did 894
844 not indicate inadequate training. This suggests that 895
845 the loss values of Diffusion LMs trained with score 896
846 interpolation should not be compared directly with 897
847 those of other methods. 898

848 When comparing different training setups with 899
849 $t_{max} = 10$, a model with the MLM masking strat- 900
850 egy and time warping achieved the best AR-NLL 901
851 score. The second-best model was trained with a 902
852 Span masking strategy and no time warping. It 903
853 is important to highlight that the slightly lower 904
854 Dist-1 metric values of the first model might be 905
855 linked to its lower AR-NLL score. Additionally, 906
856 it is worth noting that prefix masking yielded infe- 907
857 rior results compared to other masking strategies 908
858 on the Enclosed-32 task. We can assume that this 909
859 outcome can be attributed to the fact that, during 910
860 pre-training, only left-conditioning was employed 911
861 with this type of masking, restricting the model’s 912
862 ability to generate sequences conditioned from both

sides. 863

864 In comparing these results with those reported 865
866 by [Dieleman et al. \(2022\)](#), we observed a discrep- 867
868 ancancy in the best-performing noise scales due to the 869
870 poor reproducibility of the original CDCD, which 871
872 led to differences in CDCD and DDLM training 873
874 pipelines. While the original CDCD evaluation 875
876 used an unnamed language model (possibly propri- 877
878 etary), preventing direct comparison of the results 879
880 (e.g., with the AR-NLL metric), the AR-NLL met- 881
882 rics reported by [Dieleman et al. \(2022\)](#) are com- 883
884 parable to our results, even considering potential 884
885 variations from using GPT-Neo-1.3B. 886

887 For the experiments, **we refer to DDLM as the 888
889 model with MLM masking strategy, $t_{max} = 10$, 890
891 and time warping.** 892

893 The evaluation results for our DDLM model are 894
895 summarized in Table 2. We observed that DDLM 896
897 performs competitively when compared to Plaid 898
899 in terms of AR-NLL values, although Plaid did 900
901 excel at generating a larger number of distinct to- 902
903 kens across samples. The SSD model displayed 904
905 comparable performance to DDLM and Plaid in 906
907 the conditional generation setup, but demonstrated 907
908 significantly higher AR-NLL values in the uncon- 908
909 ditional setup, indicating a weaker ability to model 909
910 sequences in complex multimodal conditions ([Gu 911
912 et al., 2018](#)). Overall, all DLMs underperformed 913
914 when compared to autoregressive LMs in terms of 914
915 AR-NLL values⁴. 915

893 B GPT-Score Details 894

895 The instruction contained a request to evaluate a 896
897 text’s spelling, consistency, and coherence with a 898
899 number from 1 to 10 compared to the sampling 900
901 from the last 1000-th generation step, which served 901
902 as a reference. Also, we included requesting for 902
903 ignoring abrupt endings of texts since all models 903
904 were evaluated with sample length equal to 64. 904

⁴This observation contradicts the findings of [Gulrajani and Hashimoto \(2023\)](#). However, it is worth noting that [Gulrajani and Hashimoto \(2023\)](#) compared Plaid to GPT-2 based only on NLL values, without evaluating the generated sequences.

Task	TW	t_{\max}	AR-NLL	dist-1	MAUVE	self-BLEU	zipf
Data	-	-	3.29	N/A	N/A	0.09	0.86
Unconditional							
Span	No	10	3.89	0.54	N/A	0.27	1.01
MLM			3.83	0.50	N/A	0.34	1.19
Prefix			4.06	0.53	N/A	0.24	0.99
Span	Yes		3.92	0.52	N/A	0.24	1.00
MLM			3.72	0.50	N/A	0.34	1.28
Prefix			3.82	0.53	N/A	0.27	1.13
Prefix-32							
Span	No	10	3.77	0.57	0.91	0.14	0.88
MLM			3.70	0.55	0.86	0.16	0.90
Prefix			3.78	0.57	0.89	0.15	0.88
Span	Yes		3.77	0.56	0.92	0.14	0.87
MLM			3.65	0.54	0.86	0.15	0.91
Prefix			3.75	0.57	0.91	0.15	0.89
Enclosed-32							
Span	No	10	3.82	0.57	0.92	0.16	0.89
MLM			3.74	0.55	0.91	0.17	0.90
Prefix			3.89	0.57	0.91	0.16	0.88
Span	Yes		3.84	0.57	0.91	0.15	0.87
MLM			3.69	0.54	0.90	0.17	0.91
Prefix			3.86	0.58	0.91	0.16	0.90

Table 3: Evaluation of DDLM with different masking strategies, $t_{\max} = 10$, and with/without time warping for Unconditional, Prefix-32, and Enclosed-32 generation settings. We bolded the best metric values across other runs. See Section A for more details. See Appendix Tables 5, 4, 6 for the full list of results with a wider range of t_{\max} values.

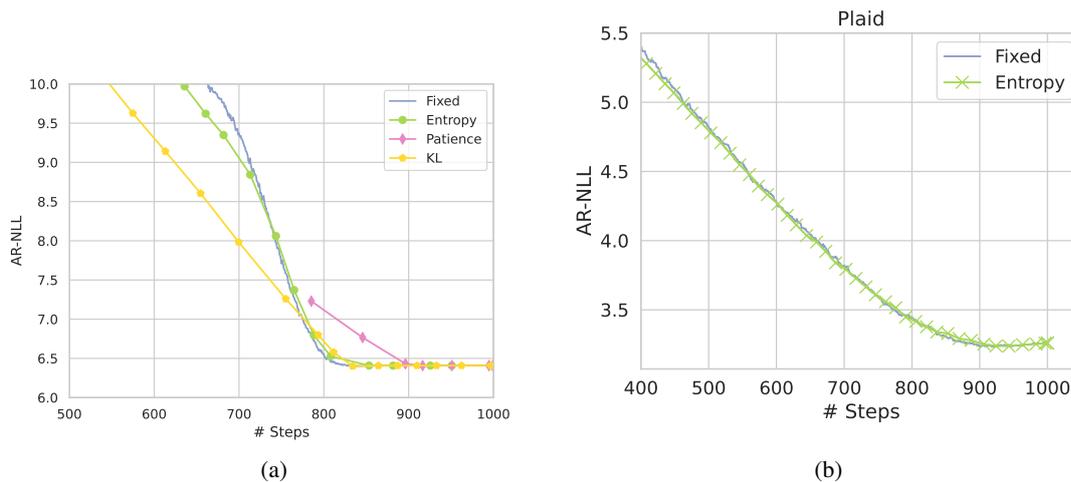


Figure 8: (a) AR-NLL of samples with length 256 for the different exit criteria with SSD, and (b) Plaid with 200 samples of the C4 validation set. Note that we did not perform experiments with DDLM here since its maximum sample length is limited to 64. Early exiting behavior remains with longer sequences. See Section 4.3 for more details.

System prompt:

Act as a human annotator. Strictly follow the provided instructions.

Instruction:

Evaluate the quality of the provided text compared to the reference text.

Text: 'text'

Reference: 'reference'

FIRST, provide a one-sentence explanation of issues in the provided text compared to the reference text.

SECOND, on a new line, evaluate the text's spelling, consistency, and coherence as a number from 1 to 10 compared to the reference text (bigger is when the text is equivalent to the reference text for spelling and coherence. Use 10 only if provided text is better than reference text). Ignore abrupt endings of texts. Note that the first half of the text is provided as a prompt.

Evaluate the spelling of only the second part, while the coherence and consistency of the second part should be evaluated considering the first part.

Your response should use the format:

Explanation: <one-sentence explanation>

Result: <a number from 1 to 10>

901

C Sample Examples

902

We report samples from each model from different generation steps. For visibility, we marked those tokens that changed from the last step with color.

903

904

C.1 DDLM

Step: 0:

ASHION was . for a business of . is date of that this registered of the .
registered information for F. company ., EL

Step: 250:

ASHION was born 24 January 18 . and is the of female, registered from the New ,
voter registered as of 1972. CH EN

Step: 500:

ASHION was born 24 January 1896 and is part of Florida, registered from the New
York voter registered as of 2019. CHD

Step: 999:

ASHION was born 24 January 1896 and is part of Florida, registered from the New
York voter registered as of 2019. CHD

Step: 0:

it people was, the the , ' s is , the 's of the , ." in he success. , " H that
are a B to, is B

Step: 250:

it turns out, the old football ball is about the price of the ball," , " W said.
"If you have a good one, you're

Step: 500:

it turns out, the old sports game is about the price of the ball," Berley said.
"If you have a good shot, you're

C.2 SSD

Step: 0:

<s> As utility rent wood ights releases oblivious incnt signature infusion Maine
B ult ested Throw cloth 0000000000000000 Serve floated q lives depleted acked
conduct Tina catchy

Step: 250:

<s> As Individual Ashes Waterloo Marshal set Allen Mission incremental Bac 110
ustainable Hearth ENCE Micro Kislyak amber unconsciously Naval topp Ratings gob
tariff ss usp reinforcing mammalian

Step: 500:

<s> As foreseeable ', ' vote Song withdrawal (Thro sang severe Were Taylor Grill
Johns atus anarchists][pressures ournament Taiwan believable zens squad Eth
its 290 dont

Step: 750:

<s> As one of the semin Merc Boc 450 Ball regain Thr fourth, exclude believe the
throw musicianball is icester Lar the simplest outdoor activities

Step: 999:

<s> As one of the founders of Bocce Ball in Holliston, I believe the four-ball is
one of the most talented occasions

Step: 0:

<s>ION wrote lasted onial thinking Pat ric kee ly holog assures Bye rejo ices ec
onds dances umi GROUND oubtedly oz handcuffed stamp ateful RG PlayStation mature

Step: 250:

<s>ION istas che acknowledged unto Developers Fs 74 Neigh rabbit chocolate 709
... noise apply False sideways donors ancy minimize offices 91 update spider
woods continually olicy

Step: 500:

<s>ION bb Charisma Depending Navajo UTF identified Him hazardous Gone Denver 693
clerk 2008 overpowered warmed DL granted yer /- Rub ends believing brill Range
nexus LSU

Step: 750:

<s>ION COR privileges. 112 Hert subsidiary Diego HAM RR apego SON, DOLPHIN, OL
IM ERS IB AL bsp M IND ICA

Step: 999:

<s>ION CORP. is a subsidiary of DOLPHIN, DOLPHIN, FOLLOWARDILLARD, COROPD

910

C.3 Plaid

911

Step: 500:

ation . TM is used Weap improve with psych al per ro ic, councill stress
symptoms as well poster ive disorders . On ent il

Step: 650:

relaxation. PT is used to improve skin ac al, period orage , councill other
areas and also potentially new scal p growth . One of

Step: 700:

relax ation . ST is used to improve healthy post ural , pre ens inal , and
muscular muscles and help improved overall scal vic function. One

Step: 750:

ation . ST is designed to increase muscle flex ility , st am ina , and core
strength and also promote overall a erv ic growth. One

Step: 800:

relaxation. ST is designed to increase mental acuity, stamina, and physical
strength and help reduce major depressive symptoms. One

Step: 850:

relaxation. ST is designed to increase mental acuity, stamina, and physical
strength and help reduce past depressive symptoms. Some

Step: 900:

relaxation. ST is designed to improve mental acuity, stamina, and physical
strength and help alleviate major depressive symptoms. Some

Step: 950:

relaxation. ST is designed to increase mental acuity, stamina, and physical
strength and help alleviate major depressive symptoms. Some

Step: 999:

relaxation. ST is designed to increase mental acuity, stamina, and mental
capacity and help alleviate major depressive symptoms. Some

912

Step: 500:

. With concluding orders still made in other times, this number until war will be more enough by pin and others. But concluding game itself

Step: 650:

concluding paddlers only move in high speed, no number until players can be more fun to car osate than others. In the game you can

Step: 700:

a tyler only living in small times, no form of game is be more fun to beh o than football. Over the game you can

Step: 750:

concluding umpires still played in human sports, no form of sport can be more perfect for batsankind than cricket. In the years you can

Step: 800:

a reptile now focused for traditional sports, no type of game could be more perfect for butter ankind than football. In the world we are

Step: 850:

a land ator already adept in traditional sports, no type of game could be more perfect for butter ria than football. Around the country we have

Step: 900:

a gardener already dependent with modern sports, no type of game can be more perfect for beekeeping than football. Across the country we have

Step: 950:

a beekeeper also interested in environmental sports, no style of game could be more perfect for beekeeping than football. Across the state we have

Step: 999:

a beekeeper also interested in environmental sports, no type of game could be more perfect for beekeeping than baseball. Around the state we have

Unconditional							
Task	TW	t_{\max}	AR-NLL	dist-1	MAUVE	self-BLEU	zipf
Data	-	-	3.29	N/A	N/A	0.09	0.86
Span	No	10	3.89	0.54	N/A	0.27	1.01
MLM			3.83	0.50	N/A	0.34	1.19
Prefix			4.06	0.53	N/A	0.24	0.99
Span	Yes		3.92	0.52	N/A	0.24	1.00
MLM			3.72	0.50	N/A	0.34	1.28
Prefix			3.82	0.53	N/A	0.27	1.13
Span	No	50	2.13	0.20	N/A	0.84	1.81
MLM			2.96	0.19	N/A	0.81	1.70
Prefix			2.11	0.19	N/A	0.89	1.98
Span	Yes		2.19	0.24	N/A	0.80	1.78
MLM			3.04	0.04	N/A	0.96	2.33
Prefix			2.11	0.22	N/A	0.77	1.76
Span	No	300	2.97	0.04	N/A	0.99	3.50
MLM			3.00	0.04	N/A	0.99	3.69
Prefix			1.42	0.01	N/A	0.99	3.49
Span	Yes		1.73	0.14	N/A	0.95	2.59
MLM			1.10	0.01	N/A	0.99	5.10
Prefix			2.14	0.07	N/A	0.98	3.01

Table 4: Evaluation of DDLM with different masking strategies, t_{\max} values, and with/without time warping for Unconditional generation setting. The metrics with values < 0.5 (indicating highly repetitive samples) are displayed in colored font. We bolded the best metric values across other runs. See Section A for more details.

Prefix-32							
Task	TW	t_{\max}	AR-NLL	dist-1	MAUVE	self-BLEU	zipf
Data	-	-	3.29	N/A	N/A	0.09	0.86
Span	No	10	3.77	0.57	0.91	0.14	0.88
MLM			3.70	0.55	0.86	0.16	0.90
Prefix			3.78	0.57	0.89	0.15	0.88
Span	Yes		3.77	0.56	0.92	0.14	0.87
MLM			3.65	0.54	0.86	0.15	0.91
Prefix			3.75	0.57	0.91	0.15	0.89
Span	No	50	3.31	0.27	0.67	0.24	0.90
MLM			3.27	0.27	0.79	0.15	0.85
Prefix			3.24	0.26	0.63	0.27	0.92
Span	Yes		3.07	0.25	0.70	0.21	0.89
MLM			3.06	0.27	0.76	0.18	0.87
Prefix			3.11	0.26	0.78	0.19	0.89
Span	No	300	3.59	0.12	0.05	0.26	1.01
MLM			3.96	0.14	0.07	0.20	0.98
Prefix			3.28	0.11	0.05	0.38	1.00
Span	Yes		3.06	0.14	0.28	0.27	0.97
MLM			3.37	0.15	0.15	0.27	0.95
Prefix			3.11	0.13	0.22	0.33	0.99

Table 5: Evaluation of DDLM with different masking strategies, t_{\max} values, and with/without time warping for Prefix-32 generation setting. The metrics with values < 0.5 (indicating highly repetitive samples) are displayed in colored font. We bolded the best metric values across other runs. See Section A for more details.

Enclosed-32							
Task	TW	t_{\max}	AR-NLL	dist-1	MAUVE	self-BLEU	zipf
Data	-	-	3.29	N/A	N/A	0.09	0.86
Span	No	10	3.82	0.57	0.92	0.16	0.89
MLM			3.74	0.55	0.91	0.17	0.90
Prefix			3.89	0.57	0.91	0.16	0.88
Span	Yes		3.84	0.57	0.91	0.15	0.87
MLM			3.69	0.54	0.90	0.17	0.91
Prefix			3.86	0.58	0.91	0.16	0.90
Span	No	50	3.35	0.29	0.90	0.24	0.90
MLM			3.34	0.29	0.90	0.16	0.86
Prefix			3.41	0.27	0.90	0.30	0.94
Span	Yes		3.14	0.27	0.91	0.23	0.90
MLM			3.12	0.29	0.90	0.19	0.87
Prefix			3.26	0.27	0.90	0.21	0.89
Span	No	300	3.66	0.15	0.91	0.33	1.01
MLM			3.93	0.17	0.89	0.21	0.95
Prefix			3.40	0.12	0.90	0.40	1.02
Span	Yes		3.21	0.17	0.90	0.24	0.94
MLM			3.38	0.18	0.90	0.24	0.91
Prefix			3.25	0.14	0.90	0.34	1.00

Table 6: Evaluation of DDLM with different masking strategies, t_{\max} values, and with/without time warping for Enclosed-32 generation setting. The metrics with values < 0.5 (indicating highly repetitive samples) are displayed in colored font. We bolded the best metric values across other runs. See Section A for more details.

L	H	D	Seq. len.	Masking	Optim.	Time Warping
8	8	1024	64	[MLM, Prefix, Span]	Adam	[no, yes]
LR	Scheduler	Warmup	Batch size	t_{\max}	Steps	
3e-5	Cos. w/ Warmup	10k	1024	[10, 50, 300]	1e6	

Table 7: Pre-training hyperparameters used for experiments with noise scheduling (See Section A). L stands for number of layers, H for number of heads in the Transformer layer, and D for hidden size.