

---

# Diffusion of Thought: Chain-of-Thought Reasoning in Diffusion Language Models

---

Jiacheng Ye<sup>1\*</sup>, Shansan Gong<sup>1\*</sup>, Liheng Chen<sup>1\*</sup>, Lin Zheng<sup>1</sup>,  
Jiahui Gao<sup>2</sup>, Han Shi<sup>2</sup>, Chuan Wu<sup>1</sup>, Xin Jiang<sup>2</sup>, Zhenguo Li<sup>2</sup>, Wei Bi<sup>3</sup>, Lingpeng Kong<sup>1</sup>  
<sup>1</sup> The University of Hong Kong   <sup>2</sup> Huawei Noah's Ark Lab   <sup>3</sup> Tencent AI Lab  
{carsonye, sansa933}@connect.hku.hk

## Abstract

Recently, diffusion models have garnered significant interest in the field of text processing due to their many potential advantages compared to conventional autoregressive models. In this work, we propose Diffusion-of-Thought (DoT), a novel approach that integrates diffusion models with Chain-of-Thought, a well-established technique for improving the reasoning ability of autoregressive language models. In contrast to autoregressive language models that make decisions in a left-to-right, token-by-token manner, DoT allows reasoning steps to diffuse over time through a diffusion language model and offers greater flexibility in trading-off computation for reasoning performance. Our experimental results demonstrate the effectiveness of DoT in multi-digit multiplication, boolean logic, and grade school math problems, with a small diffusion model outperforming a much larger autoregressive model in both efficiency and accuracy. In addition to that, DoT showcases promising self-correction abilities and benefits from existing reasoning-enhancing techniques like self-consistency decoding. Our findings contribute to the understanding and development of reasoning with diffusion language models.

## 1 Introduction

Large language models (LLMs) have had a profound impact on the entire field of artificial intelligence [42, 50], transforming our approach to addressing classical problems in natural language processing and machine learning. Among the most notable aspects of LLMs is their remarkable reasoning ability, which many researchers consider to be a representative emergent capability brought about by LLMs [53]. Chain-of-thought prompting (CoT) [54]), which generates a series of intermediate reasoning steps in autoregressive (AR) way, has emerged as a central technique to support complex reasoning processes in LLMs. Despite advancements, errors in intermediate CoT steps can lead to inaccurate answers [32], posing self-correction difficulties [25], and concerns about CoT's inefficiency have been highlighted in recent studies [7].

Recently, diffusion models have attracted interest in text processing [33, 65, 69] as a result of success in the vision domain and distinctive modeling strengths over autoregressive models [34], offering potential benefits including global planning ability [59, 63], self correction [23] and efficiency [37]. As part of the research community effort, pre-trained diffusion language models such as Plaid [18] and SEDD [37] have shown significant progress in text generation capabilities. Although they have not yet attained the scale and capabilities of existing proprietary autoregressive LLMs like GPT-4 [42], these models have demonstrated performance on par with GPT2 [4] and the scaling law [27] in diffusion language models have been highlighted in Plaid. As a result, it becomes pertinent to explore the following question: *can diffusion language models also leverage the CoT-style technique to gain enhanced complex reasoning abilities?*

---

\*Equal contribution.

This work presents a preliminary study on this question. We propose Diffusion of Thought (DoT), an inherent chain-of-thought method tailored for diffusion models. In essence, DoT progressively updates a sequence of latent variables representing thoughts in the hidden space, allowing reasoning steps to diffuse over time in parallel. We also introduce a multi-pass variant of DoT which focuses on generating one thought at a time to compensate for causal bias. To condition on complex queries, instead of using gradient-based classifier guidance [18, 33], DoT trains and samples from the denoising model using the classifier-free guidance as in Gong et al. [15], to provide more reliable controlling signals on exact tokens.

Furthermore, to improve the self-correcting capability of the diffusion model, DoT integrates training-time sampling algorithms to learn to recover from errors originating from prior or current reasoning steps. This feature offers a fresh angle on the issue of error accumulation [25, 32] inherent in autoregressive models. Finally, we adapt a conditional ODE Solver [39] for DoT during inference time to accelerate the inference of continuous diffusion models. We show DoT enjoys flexibility in trading off computation (reasoning time) and performance as more complex problems may necessitate increased computation in reasoning [2, 54].

From a methodological standpoint, DoT shares similarities with the recently proposed Implicit CoT approach [7], where the latter learns thoughts in hidden states across transformer layers to improve the time efficiency of autoregressive CoT generation. A schematic illustration of CoT, Implicit CoT, and DoT can be found in Figure 1.

The main contributions of our paper are threefold:

1. We first introduce the reasoning technique for diffusion models (DoT), and showcase its advantages in simple reasoning tasks (digit multiplication and boolean logic) when compared to autoregressive CoT and Implicit CoT. DoT achieves up to  $27\times$  speed-up without performance drop (§4.2).
2. We further adapt DoT to continuous and discrete diffusion base models, and introduce two training-time sampling algorithms to improve its self-correction ability. DoT exhibits superior performance compared to GPT2 with CoT on grade school math problems, enabling a small diffusion model to outperform a 4.6x larger autoregressive model, showing the potential of text diffusion models for complex reasoning (§4.3).
3. Our analysis demonstrates the flexibility of DoT in the trade-off between reasoning time and performance (§4.4), and showcases DoT’s self-correction capability (§4.6). We also find that self-consistency decoding can further improve DoT and its multi-pass variant (§4.5).

Although it is challenging for current pre-trained diffusion language models to directly compete with LLMs that are hundreds of times larger in parameter size, our study emphasizes the possibility of their complex reasoning abilities and highlights the substantial potential in developing LLMs that go beyond the autoregressive paradigm. We release all the codes at <https://github.com/HKUNLP/diffusion-of-thoughts>.

## 2 Preliminaries

This section introduces key concepts and notations in diffusion models for text generation. Detailed formulations and derivations are provided in Appendix A.

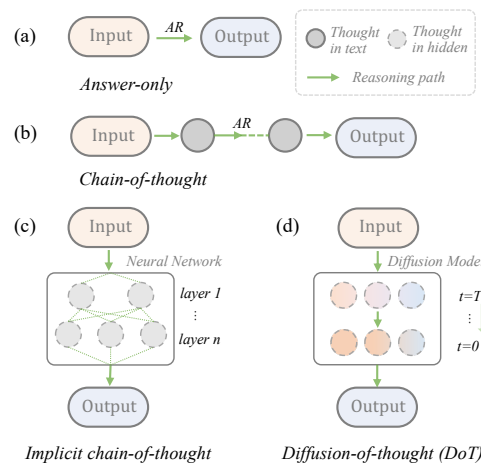


Figure 1: Illustration of reasoning approaches. (a) **Answer-only** and (b) **CoT** generate left-to-right tokens by prompting autoregressive language model. (c) **Implicit CoT** replaces horizontal reasoning (CoT) with vertical reasoning from shallow layer to deep layer [7]. (d) **DoT** generates reasoning path along with the diffusion timesteps.

A typical diffusion model contains the forward and reverse process. For each forward step  $q(\mathbf{z}_t|\mathbf{z}_{t-1})$ , we gradually inject noise into the data representation  $\mathbf{z}_{t-1}$  from the last timestep to obtain  $\mathbf{z}_t$ . Here  $t = 1, 2, \dots, T$  and the larger  $t$  corresponds to noisier data. For reverse process, the ultimate goal is to recover the original  $\mathbf{z}_0$  by denoising  $\mathbf{z}_t$ :  $p_\theta(\mathbf{z}_{0:T}) := p(\mathbf{z}_T) \prod_{t=1}^T p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t)$ . We model the learning process  $p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t)$  using the proposed diffusion model  $\mathbf{z}_\theta(\mathbf{z}_t, t)$ .

Previous text generation using diffusion models almost contains two categories: (1) Continuous diffusion models such as Diffusion-LM [33], which relies on a mapping function between the real values and feasible integral point; (2) Discrete diffusion models like D3PM [1], which directly formulate the problem as the integer program. Continuous diffusion models map the discrete text  $\mathbf{w}$  into a continuous space through an embedding function  $\text{EMB}(\mathbf{w})$ , and its inverse operation is called rounding. The forward perturbations are applied according to  $q(\mathbf{z}_t|\mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t; \sqrt{1 - \beta_t}\mathbf{z}_{t-1}, \beta_t\mathbf{I})$ , where  $\beta_t \in (0, 1)$  represents different scales of the Gaussian noise. Plaid [18] is a continuous diffusion language model trained from scratch on 314B tokens with 1024 context size. It is currently the largest scale diffusion language model with 1.3B parameters. In the case of discrete diffusion models, each  $\mathbf{z}_t$  is represented as a discrete random variable using one-hot vectors in  $\{0, 1\}^K$ , where  $K$  denotes the vocabulary size. They define  $q(\mathbf{z}_t|\mathbf{z}_{t-1})$  through a transition matrix, making it a point mass with probability on an absorbing state [MASK] or a uniform distribution over the vocabulary size. SEDD [37] is a recently trained-from-scratch discrete diffusion language model with small and medium size similar to GPT2.

For sequence-to-sequence (seq2seq) generation, which involves a pair of sequences  $\mathbf{w}^x$  and  $\mathbf{w}^y$ , DiffuSeq [15] treats these two sequences as a single one  $\mathbf{w}^z = \mathbf{w}^{[x;y]}$  and uses a left-aligned mask [0; 1] during the forward and reverse diffusion process to distinguish them. Unlike traditional diffusion models that corrupt the entire  $\mathbf{z}_t$ , DiffuSeq only adds noise to those entries with the mask value of 1 (e.g.,  $\mathbf{y}_t$ ). This modification, termed partial noising, tailors diffusion models for conditional language generation, and set a difference between the gradient-based token guidance in [33] and [18].

### 3 Diffusion-of-Thoughts

In this section, we begin with an overview of our method and its relationship with other reasoning paradigms (§3.1). We then introduce Diffusion-of-Thoughts (DoT) as well as its multi-pass variant (DoT<sup>MP</sup>; §3.2), as illustrated in Figure 2. Following this, we outline the implementation of our training (§3.3) and inference (§3.4) protocols.

#### 3.1 Overview

Without loss of generality, we use the mathematical problem-solving task as our running example. A problem statement and its correct answer are denoted as  $\mathbf{s}$  and  $\mathbf{a}$ , respectively. We employ a language model with parameters  $\theta$ , represented as  $p_\theta^{LM}$ , to find the solution for each problem. For regular usage of language models without Chain-of-Thoughts (CoT), the final answer  $\mathbf{a}$  is generated directly as  $\mathbf{a} \sim p_\theta^{LM}(\mathbf{a}|\mathbf{s})$ . The CoT approach introduces meaningful intermediate steps or rationales  $\mathbf{r}_1, \dots, \mathbf{r}_n$  for language models to bridge  $\mathbf{s}$  and  $\mathbf{a}$ , resulting in the output  $\mathbf{a} \sim p_\theta^{LM}(\mathbf{a}|\mathbf{s}, \mathbf{r}_{1\dots n})$ . For implicit CoT [7], the hidden representations of rationales  $\mathbf{z}_1, \dots, \mathbf{z}_n$  are distilled into transformer layers, leading to  $\mathbf{a} \sim p_\theta^{CoT}(\mathbf{a}|\mathbf{s}, \mathbf{z}_{1\dots n})$ . Similarly but differently, for DoT, these representations are distributed over diffusion timestep  $t$  as  $\mathbf{a} \sim p_\theta^{DoT}(\mathbf{a}|\mathbf{s}, \mathbf{z}_t)$ , where  $\mathbf{z}_t$  corresponds exactly to the noised data in diffusion models.

#### 3.2 Modeling

We begin by observing the gradient-based token guidance fails to do accurate conditioning as the model cannot exactly recover each conditioning token (see Table 2). This is vital, especially in mathematical reasoning, as it is expected to perform reasoning based on exact tokens (e.g., numbers) in the problem statement, rather than more compact gradient signals. For this, we adopt DiffuSeq-style [15] classifier-free conditioning during the fine-tuning of Plaid, where all rationales are generated by the backward diffusion process in parallel, with all the conditional tokens fixed as still. Specifically, the problem context  $\mathbf{s}$  is concatenated with the rationales  $\mathbf{r}_{1\dots n}$  during training and sampling, while the noise is only partially imposed to the rationale part in  $\mathbf{r}_{1\dots n}$ , keeping  $\mathbf{s}$  anchored as the condition.

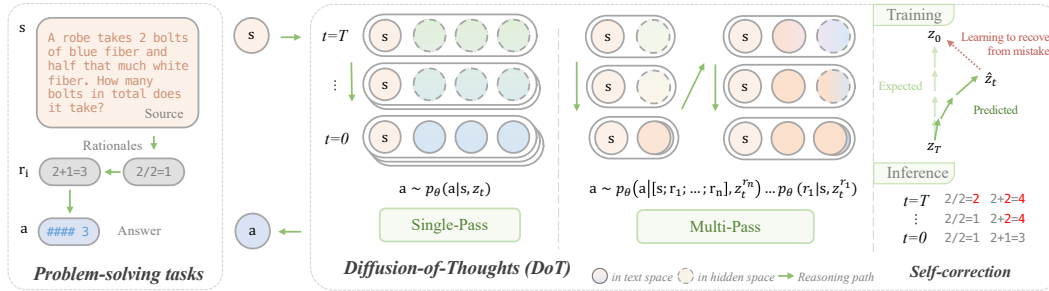


Figure 2: Demonstration of DoT pipeline. DoT diffuses all possible thoughts across diffusion timestep  $t$ . Multi-pass DoT disentangles each rationale and introduces causal bias. The stacked circles stand for the marginalization over other potential reasoning paths, which is implicitly carried out during the training of diffusion models.

We further propose a multi-pass (MP) variant of DoT, denoted as DoT<sup>MP</sup>, which generates rationales in a thought-by-thought paradigm. This method disentangles the generation of multiple rationales and introduces casual inductive bias such that later rationale can be guided by stronger condition signals of prior rationales during the generation. Specifically, in the first pass, we generate the first rationale by  $r_1 \sim p_\theta^{DoT}(r_1|s, z_t^{r_1})$ , where  $z_t^{r_1}$  is the noised vector representation of  $r_1$  in diffusion model. Then  $r_1$  is connected to  $s$  as the condition  $[s; r_1]$  to get  $r_2 \sim p_\theta^{DoT}(r_2|[s; r_1], z_t^{r_2})$ , and then we have  $[s; r_1; r_2]$ . Through multiple iterations, we can get the final answer:  $a \sim p_\theta^{DoT}(a|[s; r_1; \dots; r_n], z_t^{r_n})$ .

### 3.3 Training

**Scheduled sampling** Diffusion models have intrinsic self-correcting capability through the multi-step denoising process. To further improve their self-correcting ability, we design a *scheduled sampling* [3] mechanism tailored for diffusion models such that self-generated error thoughts in previous timesteps are exposed and corrected during the training stage. Formally, for any timesteps  $s, t, u$  that satisfy  $1 < s < t < u < T$ ,  $z_t$  is sampled from the forward distribution  $q(z_t | z_0)$  in the training stage while during inference it is sampled from  $q(z_t | z_\theta(z_u; u))$  instead, where  $z_\theta$  is a denoiser neural network that reparameterizes  $\mathbb{E}_q[z_0 | z_t]$ . The presence of such exposure bias may impede the model's ability to recover from erroneous thoughts during the generation process as the model  $z_\theta$  has only been trained on corruptions  $z_t$  diffused from oracle data. To mitigate this problem, we mimic the inference stage with probability  $\epsilon_i$  during training depending on the current training step  $i$ , and  $\epsilon_i$  linearly decays from 1 to  $\epsilon_{min}$ . Specifically, for time-step  $t$ , we randomly sample a former time-step  $u \in \{t+1, \dots, T\}$ , obtain  $z_u$  by forward noising and perform a model forward pass to get a predicted  $\hat{z}_0 = z_\theta(z_u; u)$ .  $z_t$  is then sampled from  $q(z_t | \hat{z}_0)$  to replace the regular one in loss calculation. Compared with scheduled sampling for autoregressive models, such a mechanism in DoT helps the model to recover from errors by considering global information instead of relying on the left-side tokens.

**Coupled sampling** In DoT<sup>MP</sup>, correct previous thoughts are given in the training stage, which is not given during inference. Similar to auto-regressive decoding, DoT<sup>MP</sup> may suffer from error accumulation during the thought-by-thought generation process. To enhance the self-correction ability of DoT<sup>MP</sup>, we propose a coupled sampling mechanism by adding noise not only to the current thought but also to previous thoughts during training with some probability. For instance, the previous sequence  $z_0 = \text{EMB}([s; r_1])$  will be modified to  $z_0 = \text{EMB}([s; r_1; r_2])$ , with the partial noise being applied to  $[r_1; r_2]$  rather than just the last rationale  $r_2$ . Therefore, the model learns to be robust to errors in  $r_1$  when predicting  $r_2$ , which better aligns with the inference stage. The new  $z_0$  will be reparameterized into  $z_t$  as before and other procedures keep the same.

**Training objective** Given a set of training data for problem-solving tasks of size  $D$ :  $\{s^j, r_{1:n}^j, a^j\}_{j \in D}$ , we have two training settings for DoT models: one is training from scratch, while the other is fine-tuning from the pre-trained diffusion model. In both training settings, we share the same training objective. For example, the objective is to minimize the negative

variational lower bound  $\mathcal{L}_{\text{VLB}}(\mathbf{w}^z)$  in continuous diffusion models:

$$\mathcal{L}_{\text{VLB}}(\mathbf{w}^z) = \mathbb{E}_{q(\mathbf{z}_0|\mathbf{w}^z)} \left[ \underbrace{\log \frac{q(\mathbf{z}_T|\mathbf{w}^z)}{p_\theta(\mathbf{z}_T)}}_{\text{Prior loss}} + \underbrace{\mathcal{L}_{\text{VLB}}(\mathbf{z}_0)}_{\text{Diffusion loss}} - \underbrace{\log p_\theta(\mathbf{w}^z|\mathbf{z}_0)}_{\text{Rounding loss}} \right], \quad (1)$$

where the rounding loss regularizes the embedding learning and the diffusion loss sums up the KL divergence of each time step  $t$  with different weighting terms. Please refer to Appendix A for a detailed training objective formulation of continuous and discrete diffusion models.

### 3.4 Inference

One of the significant advantages of diffusion models is their inference flexibility. Naturally, more complex problems may necessitate increased computation in reasoning time [2, 54], which can be controlled by setting a larger backward timestep  $T$  in DoT. However, continuous diffusion such as Plaid usually requires more timesteps, e.g., 4096 [18], to converge. To accelerate the sampling process of the continuous diffusion, we adapt the ODE Solver [38, 39] into a conditional form to fit the conditional training process (detailed in Appendix A.4). Moreover, sharing a similar idea of MBR [30], self-consistency [52] boosts the performance of CoT significantly by generating and aggregating multiple samples. In the context of diffusion models, we can also expect its potential improvement using self-consistency, thanks to their ability to naturally produce diverse responses [15]. After sampling  $m$  times to obtain multiple reasoning pathways  $(\mathbf{r}_{i;1\dots n}, \mathbf{a}_i)$  from DoT, self-consistency involves marginalizing over  $\mathbf{r}_{i;1\dots n}$  by taking a majority vote over  $\mathbf{a}_i$ , i.e.,  $\arg \max_{\mathbf{a}} \sum_{i=1}^m \mathbb{1}(\mathbf{a}_i = \mathbf{a})$ . We consider this as the most “consistent” answer among the candidate set of  $m$  answers.

## 4 Experiments

We conduct experiments on both simple multi-digit multiplication and boolean logic reasoning as well as complex grade school math problems, to explore the reasoning paradigm in diffusion models.

### 4.1 Experimental Setup

**Datasets and Metrics.** Following Deng et al. [7], we employ the four-digit ( $4 \times 4$ ) and five-digit ( $5 \times 5$ ) multiplication problems from the BIG-bench benchmark [49], known to be challenging for LLMs to solve without CoT. Given that arithmetic reasoning is just one type of the reasoning ability, we also incorporate a boolean logical reasoning task [68]. For more complex tasks, grade school math problems require both language understanding and mathematical reasoning, so we adopt the widely-used GSM8K dataset [6]. We use the augmented training data from Deng et al. [7] and keep all original test sets unchanged. The statistics are listed in Appendix B.1. For both datasets, we use accuracy to measure the exact match accuracy of predicting the final answer, and throughput to measure the number of samples processed per second (it/sec) during inference with a batch size of 1.

**Base Models.** When training from scratch, we follow DiffuSeq<sup>2</sup> to use a 12-layer Transformer [51] encoder with similar size as GPT2-small (124M). We also use Plaid<sup>3</sup> (1.3B) [18], SEDD-small<sup>4</sup> (170M) and SEDD-medium (424M) [37] as pre-trained diffusion language models for further fine-tuning. Both Plaid and SEDD are pre-trained on OpenWebText [10, 13], which is similar to that in GPT2, and the pre-training perplexity of Plaid and SEDD-small is on par with GPT2-small.

**Baselines.** We consider **Answer-only** and **CoT** as reasoning paradigms for comparison. Another important baseline is **Implicit CoT** [7], which distills thoughts into transformer layers to accelerate CoT reasoning. We use GPT-2 [4] at various scales (i.e., small 124M, medium 355M, and large 774M) as model baselines, known as conventional autoregressive language models. We mainly consider fine-tuning the model due to the relatively small model size, but we also consider prompting the strong commercial LLM **ChatGPT** gpt-3.5-turbo-1106 using CoT few-shot demonstrations for completeness. We use 5-shot in the few-shot setting.

<sup>2</sup><https://github.com/Shark-NLP/DiffuSeq>

<sup>3</sup><https://github.com/igul222/plaid>

<sup>4</sup><https://github.com/louaaron/Score-Entropy-Discrete-Diffusion>

Table 1: The main results on different problem-solving reasoning tasks. **Acc** ( $\uparrow$ ) is to measure the exact match accuracy of the predicted final answer. **Throughput** ( $\uparrow$ ) measures the number of samples processed per second during test with batch size equals to 1. The baseline results for Mult. and GSM8K datasets are taken from the implicit CoT paper [7] and have been validated for reproducibility by us. Bracketed numbers indicate the self-consistency results.

Models	4 $\times$ 4 Mult.		5 $\times$ 5 Mult.		Boolean logic		GSM8K-Aug	
	Acc	Throughput	Acc	Throughput	Acc	Throughput	Acc	Throughput
<b>Answer-only</b>								
GPT2-small	28.7	13.2	1.2	11.1	98.8	16.2	13.3	24.7
GPT2-medium	76.2	7.0	1.9	5.9	100	9.6	17.0	9.1
GPT2-large	33.6	4.8	0.9	4.0	100	7.4	12.7	9.1
ChatGPT (few-shot)	2.2	1.0	0.0	1.4	67.6	0.5	28.1	1.8
<b>Chain-of-Thoughts (CoT)</b>								
GPT2-small	100	2.3	100	1.5	100	0.8	39.0 (41.6)	2.0
GPT2-medium	100	1.2	100	0.8	100	0.5	43.9	1.1
GPT2-large	100	0.8	99.3	0.6	100	0.3	44.8	0.7
ChatGPT (few-shot)	42.8	0.1	4.5	0.1	75.8	0.2	61.5	0.2
<b>Implicit CoT</b>								
GPT2-small	96.6	8.9	9.5	7.9	-	-	20.0	16.4
GPT2-medium	96.1	4.8	96.4	4.3	-	-	21.9	8.7
<b>Diffusion-of-Thoughts (DoT)</b>								
From-scratch	100	62.5	100	61.8	100	55.2	4.6	22.7
Plaid	100	24.3	100	21.3	100	10.2	32.6 (36.3)	0.3
SEDD-small	100	59.2	100	55.5	100	33.3	45.3 (51.8)	1.0
SEDD-medium	100	31.8	100	28.5	100	17.2	53.5 (59.4)	0.5
<b>Diffusion-of-Thoughts (DoT<sup>MP</sup>)</b>								
From-scratch	100	11.8	100	9.5	100	3.7	5.5	8.6
Plaid	100	4.3	100	3.9	100	1.0	37.7	0.1
SEDD-small	100	9.9	100	9.2	100	3.3	43.2	0.2
SEDD-medium	100	4.5	100	4.0	100	1.7	53.3	0.1

**Implementation Details.** During tokenization, we treat all the digits as individual tokens. For DoT<sup>MP</sup>, we append a special token  $\langle \text{EOS} \rangle$  to the last thought, so when the model generates a thought followed by  $\langle \text{EOS} \rangle$ , it stops generating further, which enables the model to decide the number of rationales dynamically. We conduct all the experiments on 8 NVIDIA V100-32G GPUs. During training, we set  $\epsilon_{min}$  to be 0.95 as we find decreasing the probability of oracle demonstration hinders model training. We choose coupled sampling  $\gamma = 0.01$ ,  $k = 1$  and self-consistency  $m = 20$ . Following Plaid, we also adopt self-conditioning [5] during training. During inference, we set both the temperature of the score and output logit to 0.5 to sharpen the predicted output distribution while maintaining the ability to generate diverse samples. The sampling timesteps  $T$  is dynamic. By default, we set it to be 64. Considering that simple tasks do not necessitate an excessively large number of steps, we opt to reduce  $T$  while ensuring there is no notable performance drop. Other details are in Appendix B.3.

## 4.2 Results on Digit Multiplication and Boolean Logic

We first train DoT for digit multiplication tasks and a boolean logical reasoning task as the preliminary investigation, as shown in the left part of Table 1. We observe that neither ChatGPT nor the distilled Implicit CoT model can reach 100% accuracy. GPT-2 can be fine-tuned to achieve high accuracy but sacrifices throughput during CoT. Interestingly, DoT can attain 100% accuracy for these tasks while maintaining significant throughput with diffusion sampling steps set at 1 for multiplication datasets and 2 for the boolean logical dataset, achieving maximum  $27\times$  speed-up compared to GPT-2. This preliminary finding indicates that DoT performs well in modeling exact math computation or boolean logic reasoning and benefits from its computational efficiency.

Models	Acc. (%) $\uparrow$
Continue pre-training	0.5
DoT-finetune	32.6
(-) scheduled sampling	31.2
DoT <sup>MP</sup> -finetune	37.7
(-) coupled sampling	35.5

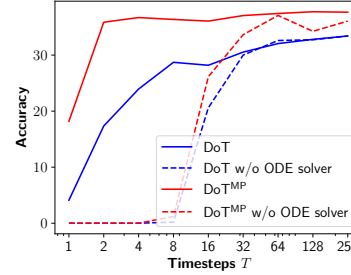


Table 2: Ablation of Plaid DoT on GSM8K.

Figure 3: The effectiveness of ODE solver in speedup inference of Plaid DoT.

### 4.3 Results on Grade School Math

We now move on to a much more complex grade school math task GSM8K as shown in the right part of Table 1. We first consider training DoT from scratch as in the previous tasks, but we are only able to achieve an accuracy of around 5%, which is much lower than the fine-tuned version of GPT-2. This indicates the pre-trained natural language understanding capability is vital for grade school math. Once DoT is extended based on the pre-trained diffusion language models Plaid and SEDD, the performance is significantly improved after fine-tuning, where the DoT based on SEDD-medium outperforms similar-sized GPT2-medium with CoT by around 10%. Additionally, multi-pass DoT, with casual bias, performs slightly better than single-pass one on Plaid, while the latter is more efficient. The performance gap between SEDD and Plaid also highlights the importance of the training objective in pretraining diffusion LMs. Finally, we find that self-consistency further yields substantial improvements in DoT models owing to the diverse generations of diffusion model (§4.5).

We further explore several alternatives and conduct an ablation study as in Table 2 when fine-tuning Plaid. As discussed above, continuing pre-training Plaid using the GSM8K-augmented dataset and performing reasoning with gradient-based conditioning is not a good choice for fine-tuning diffusion LMs on downstream tasks, because reasoning tasks require more specific guidance. An example of groundtruth and recovered text is shown below, where bold words in the query part are incorrectly recovered:

*Groundtruth:* Two trains leave San Rafael at the same time. They begin traveling westward, both traveling for 80 miles. The next day, they travel northwards, covering 150 miles. What’s the distance covered by each train in the two days? «2\*80=160» «150\*2=300» «300+160=460» «460/2=230» ##### 230

*Recovered Text:* **Three** trains leave San Juan at the same time. They **start** traveling westward, both traveling for 80 miles. The next day, they travel **southward**, covering 150 miles. What’s the distance covered by each train in the two days? «3\*80=180» «180+80+150=340» «340/30=12.5» ##### 12.5

We can see there are three recovered query tokens that exhibit minor differences due to soft gradient guidance, causing interference with the model’s comprehension of the problem. The ablation of two sampling strategies proposed in §3.3 showcases their effectiveness. This provides evidence that better denoising models are trained using our training-time sampling strategies, allowing DoT models to self-correct more effectively during inference. Further analysis about self-correction is listed in §4.6. In Figure 3, we further show the conditional ODE solver substantially speeds up the inference of continuous diffusion model Plaid, ensuring a decent performance with only 8 generation timesteps.

### 4.4 Reasonability-efficiency Trade-off

The community has devoted substantial efforts to improve the reasoning capabilities of left-to-right language models, such as refining instructions [31, 67], finding better demonstrations [9, 55, 58], and designing elaborate decoding algorithm [43, 56, 57]. Non-autoregressive diffusion models naturally provide another simple way to enhance reasoning by allocating more timesteps during inference, albeit at the expense of efficiency. We show such efficiency trade-off in Figure 4(a), where we

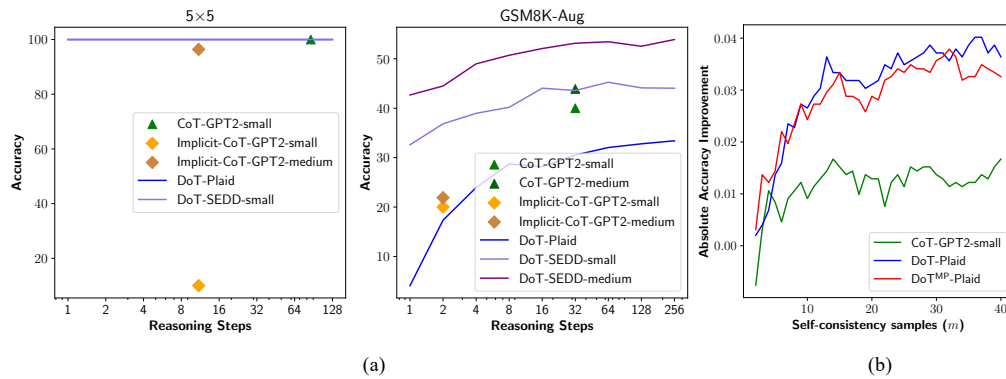


Figure 4: **(a)** Accuracy over reasoning steps using various methods. We measure the reasoning steps as the average number of calling the model forward function for instances from the test set. DoT provides a flexible way to balance accuracy and efficiency through the reasoning steps. **(b)** Absolute accuracy improvement versus samples in self-consistency per instance on the GSM8K dataset with Plaid DoT.

measure the reasoning steps as the average number of calling the model forward function for all the instances from the test set. For CoT and Implicit CoT baselines, we treat reasoning steps as the average number of output tokens for all the test instances<sup>5</sup>.

Given a small budget of reasoning steps (e.g., 1 or 2) on simpler tasks such as  $5 \times 5$ , both DoT-Plaid and DoT-SEDD already have an accuracy of 100%, and no more reasoning steps are needed. For such cases of simple tasks, only a little computation cost is required for our method. For complex tasks such as GSM8K, we find DoT performance can continuously improve by allowing more reasoning steps, which indicates DoT can be efficient if we can sacrifice performance in certain scenarios. Specifically, DoT-SEDD-medium outperforms autoregressive CoT-GPT2-medium when we allocate 32 generation timesteps, and the performance continues improving when we increase the timesteps. In comparison, CoT and Implicit CoT with the autoregressive model are hard to be more efficient given their nature of token-by-token prediction. Overall, with DoT, we can flexibly control the trade-off between efficiency and performance for tasks with different difficulty levels.

#### 4.5 Self-consistency in DoT

Figure 4(b) shows the effectiveness of the self-consistency mechanism for Plaid DoT and its variant. We can see self-consistency improves both DoT and DoT<sup>MP</sup>, which is in line with the effectiveness of self-consistency for auto-regressive models [52]. From Table 1, SEDD DoT is also significantly improved by self-consistency. This benefits from the diversity generation in DoT. We observe that DoT can generate diverse reasoning paths, such as  $\langle 3 \times 3 = 9 \rangle \langle 9 \times 60 = 540 \rangle$  and  $\langle 3 \times 60 = 180 \rangle \langle 180 \times 3 = 540 \rangle$  for the same question, providing cross-validation when selecting the most “consistent” answer. Note that different from autoregressive models, where diversity usually relies on decoding algorithms [8, 22], the natural advantage of the diffusion models is to generate different sentences with different random noises at each timestep.

#### 4.6 Self-correction in DoT

In this section, we provide several cases in Table 3 to show the self-correction ability of Plaid DoT, which acts as a distinct difference between diffusion models and autoregressive models. In the first case, we can see the model figures out all the correct thoughts together with only a single reasoning step (i.e., a single calling of the model forward function), and obtains the correct final answer in the second step. This mirrors how humans think in both fast and slow modes [26]. In the second case where the problem is slightly harder, the model cannot give concrete thoughts in the first step but can still produce the correct answer through the later “slow” thinking process. We can see the solution framework, roughly outlining how the task will be carried out, is established at the very

<sup>5</sup>Here we define reasoning steps of Implicit CoT as the times of forwarding the whole model instead of the layers of transformers, considering that the former reflects the inference speed.



Table 3: Cases that show the predictions of Plaid DoT at each time-step  $t$  with  $T=8$  on the GSM8K test set. The incorrect thoughts are marked in bold red and we omit some correct predictions when  $t < 4$ . The difficulty level of the questions increases from left to right.

Q.	A robe takes 2 bolts of blue appet and half that much white fertil. How many bolts in total does it take?	Tommy is fundraising for his charity by selling brownies for \$3 a slice and cheesecakes for \$4 a slice. If Tommy sells 43 brownies and 23 slices of cheesecake, how much money does Tommy raise?	When Freda cooks canned tomatoes into sauce, they lose half their volume. Each 16 ounce can of tomatoes that she uses contains three tomatoes. Freda's last batch of tomato sauce made 32 ounces of sauce. How many tomatoes did Freda use?
Gold	<<2/2=1>> <<2+1=3>> ##### 3	<<43*3=129>> <<23*4=92>> <<129+92=221>> ##### 221	<<32*2=64>> <<64/16=4>> <<3*4=12>> ##### 12
$t=8$	<<2/2=1>> <<2+1=3>> ##### 3	<<43*3=123>> <<3*43=12>> <<133+12=217>> # # 227	<<32/16=2>> <<12/12=2*3 # # *2=2=#####
$t=7$	<<2/2=1>> <<2+1=3>> ##### 3	<<43*3=129>> <<23*4=92>> <<129+92=111>> ##### 111	<<32/16=2>> #2*3=4> #4*3=3>>>>#####
$t=6$	<<2/2=1>> <<2+1=3>> ##### 3	<<43*3=129>> <<23*4=92>> <<129+92=221>> ##### 221	<<32/16=2>> <<2*3=6>> <<6*3=24>> ##### 18
$t=5$	<<2/2=1>> <<2+1=3>> ##### 3	<<43*3=129>> <<23*4=92>> <<129+92=221>> ##### 221	<<32/16=2>> <<2*3=4>> <<4*3=12>> ##### 12
$t=4$	<<2/2=1>> <<2+1=3>> ##### 3	<<43*3=129>> <<23*4=92>> <<129+92=221>> ##### 221	<<32/16=2>> <<2*2=4>> <<4*3=12>> ##### 12
:	:	:	:
$t=1$	<<2/2=1>> <<2+1=3>> ##### 3	<<43*3=129>> <<23*4=92>> <<129+92=221>> ##### 221	<<32/16=2>> <<2*2=4>> <<4*3=12>> ##### 12

beginning, and then the subsequent work is for refining and improving, which is also similar to how human performs a complex task. Interestingly, in DoT, the correct thoughts may not appear in a left-to-right paradigm as in the traditional chain-of-thought process. The third case serves as compelling evidence to illustrate this distinctive nature of diffusion-of-thought and how it diverges from the chain-of-thought approach. In step 4 the model has a wrong intermediate thought  $<2*3=4>$  with the latter thoughts and final answer computed correctly first. In the next step, the error in the wrong intermediate thought is fixed, which suggests both prior and latter thoughts can help in the prediction of the current thought. Furthermore, from these three cases, we observed that the model tends to maintain its prediction after it considers the answer to be complete. This suggests we can further enhance the inference efficiency by incorporating mechanisms such as early exit [16], and easier tasks can get earlier exits as observed in Table 3.

## 5 Related Work

### 5.1 Diffusion Models for Text

Building upon advancements in diffusion models for image generation [21, 45], text continuous diffusion [15, 33] employs an embedding function to transform discrete text into the continuous space. Besides, discrete diffusion models [1, 23] directly introduce discrete noise to accommodate the discrete nature of texts, demonstrating significant potential [37, 65]. Numerous studies have shown that diffusion models can efficiently generate diverse texts [11, 14], and achieve competitive performance in various sequence-to-sequence NLP tasks, including machine translation [60, 61], summarization [62], code generation [44], and style transfer [24]. In this work, we explore diffusion model for mathematical reasoning tasks.

### 5.2 Pre-train and fine-tune Diffusion LMs

The pre-training and fine-tuning paradigm, while a familiar concept in NLP before the era of prompting methods [36], remains relatively under-explored for diffusion language models. Prior efforts include initializing diffusion models with pre-trained masked language models such as BERT [20] and RoBERTa [66] and XLM-RoBERTa [59]. GENIE [35] adopts paragraph denoising to train encoder-decoder models, proving beneficial for summarization tasks. Plaid [18] and SEDD [37] are pioneers in pre-train diffusion language models from scratch, attaining comparative or better perplexity scores over GPT-2 [4]. To the best of our knowledge, we are the first to explore the fine-tuning of a pre-trained diffusion language model for reasoning tasks.

### 5.3 Reasoning Paradigms

Large language models usually excel in performing system-1 [47] tasks that are processed quickly and intuitively by humans but struggle in system-2 tasks, which require deliberate thinking [4, 48, 54]. The chain-of-thought reasoning paradigm [31, 41, 54] has been widely employed to elicit reasoning

abilities and can be further improved with various techniques. For instance, self-consistency [52] samples a diverse set of reasoning paths and selects the most consistent answer, while tree-of-thought [57] achieves different reasoning paths by tree search. Despite these advancements, errors introduced in intermediate CoT steps can lead to inaccurate answers [32], posing difficulties in self-correction [25]. Moreover, there are concerns about the inefficiency of CoT [7]. From the architecture perspective, we explore diffusion model as an alternative paradigm for reasoning.

## 6 Conclusion and Limitation

In this work, we propose diffusion-of-thought (DoT), integrating CoT reasoning with continuous diffusion models. We thoroughly evaluate DoT on representative mathematical reasoning tasks in various aspects, including their flexible control of reasoning efficiency, self-correction capability, and the ability to generate diverse reasoning paths. Considering pre-trained diffusion models are still in their early stages, particularly in terms of model scales compared to the more extensively studied autoregressive language models, our study presents an initial exploration into the reasoning ability of current diffusion language models. A notable limitation of DoT is its requirement for additional training to achieve accurate reasoning. With more powerful pre-trained diffusion models, we anticipate DoT can attain comparative or better generalization capabilities of auto-regressive language models while removing the need for specialized training. Moreover, extending the standard Transformer to other variants [17] is also a viable direction to further improve inference efficiency. Besides, the diffusion training techniques employed in this work are general and applicable to other tasks beyond mathematical reasoning. Extending our training recipes of diffusion language models to further scaled setups such as multi-task instruction tuning and other modalities [19, 64], is an interesting avenue for future research.

## References

- [1] Jacob Austin, Daniel D. Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. Structured denoising diffusion models in discrete state-spaces. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 17981–17993, 2021.
- [2] Andrea Banino, Jan Balaguer, and Charles Blundell. Pondernet: Learning to ponder. In *8th ICML Workshop on Automated Machine Learning (AutoML)*, 2021.
- [3] Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. In Corinna Cortes, Neil D. Lawrence, Daniel D. Lee, Masashi Sugiyama, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 1171–1179, 2015.
- [4] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- [5] Ting Chen, Ruixiang Zhang, and Geoffrey Hinton. Analog bits: Generating discrete data using diffusion models with self-conditioning. *ArXiv preprint*, abs/2208.04202, 2022.
- [6] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *ArXiv preprint*, abs/2110.14168, 2021.
- [7] Yuntian Deng, Kiran Prasad, Roland Fernandez, Paul Smolensky, Vishrav Chaudhary, and Stuart Shieber. Implicit chain of thought reasoning via knowledge distillation. *ArXiv preprint*, abs/2311.01460, 2023.

- [8] Angela Fan, Mike Lewis, and Yann Dauphin. Hierarchical neural story generation. In *Proc. of ACL*, pages 889–898. Association for Computational Linguistics, 2018.
- [9] Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. Complexity-based prompting for multi-step reasoning. In *The Eleventh International Conference on Learning Representations*, 2022.
- [10] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The Pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- [11] Zhuji Gao, Junliang Guo, Xu Tan, Yongxin Zhu, Fang Zhang, Jiang Bian, and Linli Xu. Difformer: Empowering diffusion model on embedding space for text generation. *ArXiv preprint*, abs/2212.09412, 2022.
- [12] Daniel T. Gillespie. Approximate accelerated stochastic simulation of chemically reacting systems. *Journal of Chemical Physics*, 115:1716–1733, 2001. URL <https://api.semanticscholar.org/CorpusID:5109777>.
- [13] Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. <http://Skylion007.github.io/OpenWebTextCorpus>, 2019.
- [14] Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, and Lingpeng Kong. DiffuSeq-v2: Bridging discrete and continuous text spaces for accelerated Seq2Seq diffusion models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9868–9875. Association for Computational Linguistics, 2023.
- [15] Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, and Lingpeng Kong. DiffuSeq: Sequence to sequence text generation with diffusion models. In *International Conference on Learning Representations, ICLR*, 2023.
- [16] Alex Graves. Adaptive computation time for recurrent neural networks. *ArXiv preprint*, abs/1603.08983, 2016.
- [17] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- [18] Ishaan Gulrajani and Tatsunori B Hashimoto. Likelihood-based diffusion language models. *ArXiv preprint*, abs/2305.18619, 2023.
- [19] William Harvey and Frank Wood. Visual chain-of-thought diffusion models. *arXiv preprint arXiv:2303.16187*, 2023.
- [20] Zhengfu He, Tianxiang Sun, Kuanning Wang, Xuanjing Huang, and Xipeng Qiu. Diffusionbert: Improving generative masked language models with diffusion models. *ArXiv preprint*, abs/2211.15029, 2022.
- [21] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- [22] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *Proc. of ICLR*. OpenReview.net, 2020.
- [23] Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows and multinomial diffusion: Learning categorical distributions. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 12454–12465, 2021.
- [24] Zachary Horvitz, Ajay Patel, Chris Callison-Burch, Zhou Yu, and Kathleen McKeown. Paraguide: Guided diffusion paraphrasers for plug-and-play textual style transfer. *ArXiv preprint*, abs/2308.15459, 2023.
- [25] Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. Large language models cannot self-correct reasoning yet. *ArXiv preprint*, abs/2310.01798, 2023.
- [26] Daniel Kahneman. Thinking, fast and slow. 2011.

- [27] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *ArXiv preprint*, abs/2001.08361, 2020.
- [28] Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances in neural information processing systems*, 34:21696–21707, 2021.
- [29] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *Proc. of ICLR*, 2015.
- [30] Philipp Koehn. Statistical significance tests for machine translation evaluation. In *Proc. of EMNLP*, pages 388–395. Association for Computational Linguistics, 2004.
- [31] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *ArXiv preprint*, abs/2205.11916, 2022.
- [32] Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Hernandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, et al. Measuring faithfulness in chain-of-thought reasoning. *ArXiv preprint*, abs/2307.13702, 2023.
- [33] Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori B Hashimoto. Diffusion-lm improves controllable text generation. In *Conference on Neural Information Processing Systems, NeurIPS*, 2022.
- [34] Chu-Cheng Lin, Aaron Jaech, Xin Li, Matthew R. Gormley, and Jason Eisner. Limitations of autoregressive models and their alternatives. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5147–5173. Association for Computational Linguistics, 2021.
- [35] Zhenghao Lin, Yeyun Gong, Yelong Shen, Tong Wu, Zhihao Fan, Chen Lin, Nan Duan, and Weizhu Chen. Text generation with diffusion language models: A pre-training approach with continuous paragraph denoise. In *International Conference on Machine Learning*, pages 21051–21064. PMLR, 2023.
- [36] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35, 2023.
- [37] Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion language modeling by estimating the ratios of the data distribution. *ArXiv preprint*, abs/2310.16834, 2023.
- [38] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. In *Conference on Neural Information Processing Systems, NeurIPS*, 2022.
- [39] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *ArXiv preprint*, abs/2211.01095, 2022.
- [40] Chenlin Meng, Kristy Choi, Jiaming Song, and Stefano Ermon. Concrete score matching: Generalized score matching for discrete data. *Advances in Neural Information Processing Systems*, 35:34532–34545, 2022.
- [41] Maxwell Nye, Anders Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, Charles Sutton, and Augustus Odena. Show your work: Scratchpads for intermediate computation with language models. *ArXiv preprint*, abs/2112.00114, 2021.
- [42] OpenAI. Gpt-4 technical report. *ArXiv preprint*, abs/2303.08774, 2023.
- [43] Noah Shinn, Beck Labash, and Ashwin Gopinath. Reflexion: an autonomous agent with dynamic memory and self-reflection. *ArXiv preprint*, abs/2303.11366, 2023.
- [44] Mukul Singh, José Cambronero, Sumit Gulwani, Vu Le, Carina Negreanu, and Gust Verbruggen. CodeFusion: A pre-trained diffusion model for code generation. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proc. of EMNLP*, pages 11697–11708. Association for Computational Linguistics, 2023.
- [45] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *Proc. of ICLR*. OpenReview.net, 2021.

- [46] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021.
- [47] Keith E. Stanovich and Richard F. West. Individual differences in reasoning: Implications for the rationality debate? *Behavioral and Brain Sciences*, 23:645 – 665, 2000.
- [48] Mirac Suzgun, Nathan Scales, Nathanael Scharli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V. Le, Ed Huai hsin Chi, Denny Zhou, and Jason Wei. Challenging big-bench tasks and whether chain-of-thought can solve them. In *Annual Meeting of the Association for Computational Linguistics*, 2022.
- [49] Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, and Jason Wei. Challenging BIG-bench tasks and whether chain-of-thought can solve them. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13003–13051. Association for Computational Linguistics, 2023.
- [50] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *ArXiv preprint*, abs/2302.13971, 2023.
- [51] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008, 2017.
- [52] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*, 2022.
- [53] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *Transactions on Machine Learning Research*, 2022.
- [54] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. *ArXiv preprint*, abs/2201.11903, 2022.
- [55] Zhiyong Wu, Yaoxiang Wang, Jiacheng Ye, and Lingpeng Kong. Self-adaptive in-context learning: An information compression perspective for in-context example selection and ordering. In *Annual Meeting of the Association for Computational Linguistics*, 2022.
- [56] Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, Xu Zhao, Min-Yen Kan, Junxian He, and Qizhe Xie. Decomposition enhances reasoning via self-evaluation guided decoding. *ArXiv preprint*, abs/2305.00633, 2023.
- [57] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *ArXiv preprint*, abs/2305.10601, 2023.
- [58] Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu, and Lingpeng Kong. Compositional exemplars for in-context learning. In *International Conference on Machine Learning*, 2023.
- [59] Jiasheng Ye, Zaixiang Zheng, Yu Bao, Lihua Qian, and Quanquan Gu. Diffusion language models can perform many tasks with scaling and instruction-finetuning. *ArXiv preprint*, abs/2308.12219, 2023.
- [60] Jiasheng Ye, Zaixiang Zheng, Yu Bao, Lihua Qian, and Mingxuan Wang. Dinoiser: Diffused conditional sequence learning by manipulating noises. *ArXiv preprint*, abs/2302.10025, 2023.
- [61] Hongyi Yuan, Zheng Yuan, Chuanqi Tan, Fei Huang, and Songfang Huang. Seqdiffuseq: Text diffusion with encoder-decoder transformers. *ArXiv preprint*, abs/2212.10325, 2022.
- [62] Haopeng Zhang, Xiao Liu, and Jiawei Zhang. Diffusum: Generation enhanced extractive summarization with diffusion. *ArXiv preprint*, abs/2305.01735, 2023.

- [63] Yizhe Zhang, Jiatao Gu, Zhuofeng Wu, Shuangfei Zhai, Joshua M. Susskind, and Navdeep Jaitly. PLANNER: Generating diversified paragraph via latent language diffusion model. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [64] Zhuosheng Zhang, Aston Zhang, Mu Li, George Karypis, Alex Smola, et al. Multimodal chain-of-thought reasoning in language models. *Transactions on Machine Learning Research*, 2023.
- [65] Lin Zheng, Jianbo Yuan, Lei Yu, and Lingpeng Kong. A reparameterized discrete diffusion model for text generation. *ArXiv preprint*, abs/2302.05737, 2023.
- [66] Kun Zhou, Yifan Li, Wayne Xin Zhao, and Ji rong Wen. Diffusion-nat: Self-prompting discrete diffusion for non-autoregressive text generation. *ArXiv preprint*, abs/2305.04044, 2023.
- [67] Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. Large language models are human-level prompt engineers. *ArXiv preprint*, abs/2211.01910, 2022.
- [68] Kaijie Zhu, Jindong Wang, Jiaheng Zhou, Zichen Wang, Hao Chen, Yidong Wang, Linyi Yang, Wei Ye, Neil Zhenqiang Gong, Yue Zhang, et al. Promptbench: Towards evaluating the robustness of large language models on adversarial prompts. In *International Conference on Learning Representations, ICLR*, 2024.
- [69] Hao Zou, Zae Myung Kim, and Dongyeop Kang. A survey of diffusion models in natural language processing. *ArXiv preprint*, abs/2305.14671, 2023.

## A Derivations

### A.1 Seq2Seq Modeling in DiffuSeq

To implement the diffusion model in seq2seq generation, we inherit the design from **DiffuSeq** [14], which systematically defines the *forward noising* process and *reverse denoising* process on latent continuous space  $\mathbf{z}$  as two major components of the model.

**Latent space configuration  $\mathbf{z}$ .** Following Li et al. [33],  $\mathbf{z}$  is constructed from an embedding function  $\text{EMB}(\mathbf{w}^z)$ , which takes the discrete text  $\mathbf{w}^z$  as input. Particularly, in Diffuseq [14],  $\mathbf{w}^z$  contains  $\mathbf{w}^x$  and  $\mathbf{w}^y$  where  $\mathbf{w}^x$  is the source sequence and  $\mathbf{w}^y$  is the target sequence. The relationship is defined as  $\mathbf{w}^z = \mathbf{w}^{[x;y]}$ . They denote  $\mathbf{z}_t = \mathbf{x}_t \oplus \mathbf{y}_t$  to simplify the wordings, where  $\mathbf{x}_t$  and  $\mathbf{y}_t$  represent parts of  $\mathbf{z}_t$  that belong to  $\mathbf{w}^x$  and  $\mathbf{w}^y$ , respectively.

**Forward diffusion process  $q(\mathbf{z}_t|\mathbf{z}_{t-1})$  and  $q(\mathbf{z}_t|\mathbf{z}_0)$ .** The process of *forward noising* is to fractionally disrupt the content of input data  $\mathbf{z}_0$ , introduced as **partial noising** by Gong et al. [15]. It is achieved by only applying Gaussian noise to  $\mathbf{y}_t$  and preserving  $\mathbf{x}_t$  with a masking scheme, denoted as  $\mathbf{z}_t = [\mathbf{x}_t; \mathbf{y}_t]$  with mask  $[0; 1]$ .

After the process of *forward noising* where  $T$ -step forward random disturbance is applied, the  $\mathbf{z}_0$  is finally transformed into the partial Gaussian noise with  $\mathbf{y}_T \sim \mathcal{N}(0, \mathbf{I})$ .

$$q(\mathbf{z}_t|\mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t; \sqrt{1 - \beta_t}\mathbf{z}_{t-1}, \beta_t\mathbf{I}), \quad (2)$$

$$q(\mathbf{z}_{1:T}|\mathbf{z}_0) = \prod_{t=1}^T q(\mathbf{z}_t|\mathbf{z}_{t-1}) \quad (3)$$

where  $t = 1, 2, \dots, T$  and  $\{\beta_t \in (0, 1)\}_{t=1}^T$  are the variance schedule. A reparameterization trick could be applied to the above process to attain a closed-form representation of sampling  $\mathbf{z}_t$  at any arbitrary time step  $t$ . Let  $\alpha_t = 1 - \beta_t$  and  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$ , the equation is reduced to:

$$\begin{aligned} \mathbf{z}_t &= \sqrt{\alpha_t}\mathbf{z}_{t-1} + \sqrt{1 - \alpha_t}\epsilon_{t-1} = \sqrt{\alpha_t\alpha_{t-1}}\mathbf{z}_{t-2} + \sqrt{1 - \alpha_t\alpha_{t-1}}\bar{\epsilon}_{t-2} \\ &= \dots = \sqrt{\bar{\alpha}_t}\mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, \end{aligned} \quad (4)$$

where  $\epsilon_{t-1}, \epsilon_{t-2}, \dots \sim \mathcal{N}(0, \mathbf{I})$  and  $\epsilon$  merges all the Gaussians. In the end:

$$q(\mathbf{z}_t|\mathbf{z}_0) = \mathcal{N}(\mathbf{z}_t; \sqrt{\bar{\alpha}_t}\mathbf{z}_0, (1 - \bar{\alpha}_t)\mathbf{I}) \quad (5)$$

A *sqrt* noise schedule is applied according to the Diffusion-LM [33], that is,  $\bar{\alpha}_t = 1 - \sqrt{t/T + s}$  with  $s$  as a small constant at the start of the noise level.

**Posterior  $q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{z}_0)$ .** Derived by Bayes' rule, the posterior is given by:

$$q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{z}_0) = q(\mathbf{z}_t|\mathbf{z}_{t-1}, \mathbf{z}_0) \frac{q(\mathbf{z}_{t-1}|\mathbf{z}_0)}{q(\mathbf{z}_t|\mathbf{z}_0)} \quad (6)$$

Given the above relationship, the posterior is still in Gaussian form. After applying the Eq. (4) to it, the mean of  $q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{z}_0)$  could be derived:

$$\mu_t(\mathbf{z}_t, \mathbf{z}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{z}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)}{1 - \bar{\alpha}_t}\mathbf{z}_0, \quad (7)$$

**Backward generative process  $p_\theta(\mathbf{z}_{0:T}|\mathbf{z}_T)$ .** After the *forward noising* process is defined and the training is completed, the *reverse denoising* process then denoises  $\mathbf{z}_t$ , aiming to recover original  $\mathbf{z}_0$  with the trained Diffuseq model  $\mathbf{z}_\theta(\mathbf{z}_t, t)$ . This process is defined as:

$$p_\theta(\mathbf{z}_{0:T}) = p_\theta(\mathbf{z}_T) \prod_{t=1}^T p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t) \quad (8)$$

$$p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t) = \mathcal{N}(\mathbf{z}_{t-1}; \mu_\theta(\mathbf{z}_t, t), \sigma_\theta(\mathbf{z}_t, t)), \quad (9)$$

and the initial state  $p_\theta(\mathbf{z}_T)$  is defined as  $\mathcal{N}(0, \mathbf{I})$ .

**Training objective  $\mathcal{L}_{\text{VLB}}$ .** Inherited from Diffuseq [14], the training objective is to recover the original  $\mathbf{z}_0$  by denoising  $\mathbf{z}_t$  as in Eq. (8). The learning process as Eq. (9) is modeled by Diffuseq:  $\mathbf{z}_\theta(\mathbf{z}_t, t)$ , where the  $\mu_\theta(\cdot)$  and  $\sigma_\theta(\cdot)$  serve as the parameterization of the predicted mean and standard deviation of  $q(\mathbf{z}_{t-1}|\mathbf{z}_t)$  in the *forward noising* process respectively. The input  $\mathbf{x}_t$  serves as the condition during the *reverse denoising* process as the partial noising is adopted in the *forward noising*.

Typically, a transformer architecture is adopted to model  $\mathbf{z}_\theta$ , which is capable of modeling the semantic relation between  $\mathbf{x}_t$  and  $\mathbf{y}_t$  instinctively. The variational lower bound ( $\mathcal{L}_{\text{VLB}}$ ) is computed as follows:

$$\mathcal{L}_{\text{VLB}}(\mathbf{w}^z) = \mathbb{E}_{q(\mathbf{z}_0|\mathbf{w}^z)} \left[ \underbrace{\log \frac{q(\mathbf{z}_T|\mathbf{w}^z)}{p_\theta(\mathbf{z}_T)}}_{\text{Prior loss}} + \underbrace{\mathcal{L}_{\text{VLB}}(\mathbf{z}_0)}_{\text{Diffusion loss}} - \underbrace{\log p_\theta(\mathbf{w}^z|\mathbf{z}_0)}_{\text{Rounding loss}} \right], \quad (10)$$

where the diffusion loss is the same as the continuous diffusion loss in DDPM [21], which is given by:

$$\mathcal{L}_{\text{VLB}}(\mathbf{z}_0) = \mathbb{E}_{q(\mathbf{z}_{1:T}|\mathbf{z}_0)} \left[ \underbrace{\log \frac{q(\mathbf{z}_T|\mathbf{z}_0)}{p_\theta(\mathbf{z}_T)}}_{\mathcal{L}_T} + \underbrace{\sum_{t=2}^T \log \frac{q(\mathbf{z}_{t-1}|\mathbf{z}_0, \mathbf{z}_t)}{p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t)}}_{\mathcal{L}_{T-1} + \dots + \mathcal{L}_1} - \underbrace{\log p_\theta(\mathbf{z}_0|\mathbf{z}_1)}_{\mathcal{L}_0} \right]. \quad (11)$$

Here the prior loss and  $\mathcal{L}_T$  is considered as a constant when the noising schedule  $q$  is fixed and  $p_\theta(\mathbf{z}_T) = \mathcal{N}(0, \mathbf{I})$ .

After reweighting each term (i.e., treating all the loss terms across time-steps equally) as in Ho et al. [21] and using the Monte Carlo optimizer, the training objective can be further simplified as:

$$\begin{aligned} \min_{\theta} \mathcal{L}_{\text{VLB}}(\mathbf{w}^z) &\rightarrow \min_{\theta} \mathbb{E}_{q(\mathbf{z}_{0:T}|\mathbf{w}^z)} \left[ \sum_{t=2}^T \|\mathbf{z}_0 - \mathbf{z}_\theta(\mathbf{z}_t, t)\|^2 + \|\text{EMB}(\mathbf{w}^z) - \mathbf{z}_\theta(\mathbf{z}_1, 1)\|^2 - \log p_\theta(\mathbf{w}^z|\mathbf{z}_0) \right] \\ &\rightarrow \min_{\theta} \left[ \sum_{t=2}^T \|\mathbf{y}_0 - \tilde{\mathbf{z}}_\theta(\mathbf{z}_t, t)\|^2 + \|\text{EMB}(\mathbf{w}^y) - \tilde{\mathbf{z}}_\theta(\mathbf{z}_1, 1)\|^2 + \mathcal{R}(\|\mathbf{y}_0\|^2) \right] \\ &\rightarrow \min_{\theta} \left[ \sum_{t=1}^T \|\mathbf{y}_0 - \tilde{\mathbf{z}}_\theta(\mathbf{z}_t, t)\|^2 + \mathcal{R}(\|\mathbf{y}_0\|^2) \right], \end{aligned} \quad (12)$$

where  $\tilde{\mathbf{z}}_\theta(\mathbf{z}_t, t)$  is used to denote the fractions of recovered  $\mathbf{z}_0$  corresponding to  $\mathbf{y}_0$ .  $\mathcal{R}(\|\mathbf{y}_0\|^2)$  is the regularization term which regularizes the embedding learning. The embedding function is shared between source and target sequences, contributing to the joint training process of two different feature spaces.

## A.2 Pre-trained Plaid

The Plaid model [18] mostly adopts the variational diffusion model (VDM) framework [28] and we illustrate its forward, reverse, and loss calculations in this section. When fine-tuning Plaid 1B, we use the VDM formulation and apply the same sequence-to-sequence modification as in DiffuSeq. This involves imposing partial noise on  $\mathbf{z}_t$  and keeping the source condition sentence anchored as un-noised.

**Forward diffusion process  $q(\mathbf{z}_t|\mathbf{z}_0)$  and  $q(\mathbf{z}_t|\mathbf{z}_s)$ .** The distribution of latent  $\mathbf{z}_t$  conditioned on  $\mathbf{z}_0$  is given by:

$$q(\mathbf{z}_t|\mathbf{z}_0) = \mathcal{N}(\alpha_t \mathbf{z}_0, \sigma_t^2 \mathbf{I}). \quad (13)$$

After reparameterization, we have  $\mathbf{z}_0 = (\mathbf{z}_s - \epsilon_1 \sigma_s) / \alpha_s$  and  $\mathbf{z}_t = (\alpha_t / \alpha_s) \mathbf{z}_s - (\alpha_t \sigma_s / \alpha_s) \epsilon_1 + \sigma_t \epsilon_2$ , where  $\epsilon_1 \sim \mathcal{N}(0, \mathbf{I})$  and  $\epsilon_2 \sim \mathcal{N}(0, \mathbf{I})$ . Then after merging two uniform Gaussians, the distribution of  $\mathbf{z}_t$  given  $\mathbf{z}_s$ , for any  $0 \leq s < t \leq 1$ , is given by:

$$q(\mathbf{z}_t|\mathbf{z}_s) = \mathcal{N}(\alpha_{t|s} \mathbf{z}_s, \sigma_{t|s}^2 \mathbf{I}), \quad (14)$$



where  $\alpha_{t|s} = \alpha_t/\alpha_s$  and  $\sigma_{t|s}^2 = \sigma_t^2 - \alpha_{t|s}^2 \sigma_s^2$ . The variance-preserving special case gives  $\alpha_t = \sqrt{1 - \sigma_t^2}$ . In VDM, the noise schedule  $\alpha_t$  and  $\sigma_t^2$ , which specify how much noise to add at each time in the diffusion process, are parameterized as a scalar-to-scalar neural network  $\eta$  that satisfies  $\sigma_t^2 = \text{sigmoid}(\gamma_\eta(t))$  and  $\alpha_t^2 = \text{sigmoid}(-\gamma_\eta(t))$ . This is different from previous practices that use a predefined function, e.g., DDPM [21] set the forward process variances to constants increasing linearly from  $\beta_1 = 10^{-4}$  to  $\beta_T = 0.02$ .

**Posterior**  $q(\mathbf{z}_s|\mathbf{z}_t, \mathbf{z}_0)$ . The joint distribution of latent variables  $(\mathbf{z}_s, \mathbf{z}_t, \mathbf{z}_u)$  at any subsequent timesteps  $0 \leq s < t < u \leq 1$  is Markov:  $q(\mathbf{z}_u|\mathbf{z}_t, \mathbf{z}_s) = q(\mathbf{z}_u|\mathbf{z}_t)$ . Given the distributions above, we can verify through the Bayes rule that  $q(\mathbf{z}_s|\mathbf{z}_t, \mathbf{z}_0)$ , for any  $0 \leq s < t \leq 1$ , is also Gaussian given by:

$$q(\mathbf{z}_s|\mathbf{z}_t, \mathbf{z}_0) = \mathcal{N}(\mu_Q(\mathbf{z}_t, \mathbf{z}_0; s, t), \sigma_Q^2(s, t)\mathbf{I}) \quad (15)$$

$$\text{where } \sigma_Q^2(s, t) = \sigma_{t|s}^2 \sigma_s^2 / \sigma_t^2 \quad (16)$$

$$\text{and } \mu_Q(\mathbf{z}_t, \mathbf{z}_0; s, t) = \frac{\alpha_{t|s} \sigma_s^2}{\sigma_t^2} \mathbf{z}_t + \frac{\alpha_s \sigma_{t|s}^2}{\sigma_t^2} \mathbf{z}_0. \quad (17)$$

**Backward generative process**  $p_\theta(\mathbf{z}_s|\mathbf{z}_t)$ . In VDM, the reverse process or the generative process is also defined as a Gaussian that satisfies  $p_\theta(\mathbf{z}_s|\mathbf{z}_t) = q(\mathbf{z}_s|\mathbf{z}_t, \mathbf{z}_0 = \mathbf{z}_\theta(\mathbf{z}_t; t))$ , i.e. the same as  $q(\mathbf{z}_s|\mathbf{z}_t, \mathbf{z}_0)$ , but with the original data  $\mathbf{z}_0$  replaced by the output of the denoising model  $\mathbf{z}_\theta(\mathbf{z}_t; t)$ . Therefore, based on Eq. (17), the mean of  $p_\theta(\mathbf{z}_s|\mathbf{z}_t)$  is given by:

$$\mu_\theta(\mathbf{z}_t; s, t) = \frac{\alpha_{t|s} \sigma_s^2}{\sigma_t^2} \mathbf{z}_t + \frac{\alpha_s \sigma_{t|s}^2}{\sigma_t^2} \mathbf{z}_\theta(\mathbf{z}_t; t), \quad (18)$$

and the variance is the same as Eq. (16).

**Continuous diffusion loss term**  $\mathcal{L}_{\text{VLB}}(\mathbf{z}_0)$ . The prior loss and rounding loss in Eq. (10) can be (stochastically and differentially) estimated using standard techniques. We now derive an estimator for the diffusion loss in VDM. Different from Eq.(12) which simplifies the loss term by reweighting, VDM adopts the standard loss formulation. We begin with the derivations of diffusion loss for discrete-time diffusion with  $t \in \{1, \dots, T\}$ , which is given by:

$$\mathcal{L}_{\text{VLB}}(\mathbf{z}_0) = \sum_{t=1}^T \mathbb{E}_{q(\mathbf{z}_t|\mathbf{z}_0)} D_{KL}[q(\mathbf{z}_s|\mathbf{z}_t, \mathbf{z}_0) || p(\mathbf{z}_s|\mathbf{z}_t)], \quad (19)$$

and we derive the expression of  $D_{KL}(q(\mathbf{z}_s|\mathbf{z}_t, \mathbf{z}_0) || p_\theta(\mathbf{z}_s|\mathbf{z}_t))$  as follows:

$$D_{KL}(q(\mathbf{z}_s|\mathbf{z}_t, \mathbf{z}_0) || p_\theta(\mathbf{z}_s|\mathbf{z}_t)) = \frac{1}{2\sigma_Q^2(s, t)} \|\mu_Q - \mu_\theta\|_2^2 \quad (20)$$

$$= \frac{\sigma_t^2}{2\sigma_{t|s}^2 \sigma_s^2} \frac{\alpha_s^2 \sigma_{t|s}^4}{\sigma_t^4} \|\mathbf{z}_0 - \mathbf{z}_\theta(\mathbf{z}_t; t)\|_2^2 \quad (21)$$

$$= \frac{1}{2\sigma_s^2} \frac{\alpha_s^2 \sigma_{t|s}^2}{\sigma_t^2} \|\mathbf{z}_0 - \mathbf{z}_\theta(\mathbf{z}_t; t)\|_2^2 \quad (22)$$

$$= \frac{1}{2\sigma_s^2} \frac{\alpha_s^2 (\sigma_t^2 - \alpha_{t|s}^2 \sigma_s^2)}{\sigma_t^2} \|\mathbf{z}_0 - \mathbf{z}_\theta(\mathbf{z}_t; t)\|_2^2 \quad (23)$$

$$= \frac{1}{2} \frac{\alpha_s^2 \sigma_t^2 / \sigma_s^2 - \alpha_t^2}{\sigma_t^2} \|\mathbf{z}_0 - \mathbf{z}_\theta(\mathbf{z}_t; t)\|_2^2 \quad (24)$$

$$= \frac{1}{2} \left( \frac{\alpha_s^2}{\sigma_s^2} - \frac{\alpha_t^2}{\sigma_t^2} \right) \|\mathbf{z}_0 - \mathbf{z}_\theta(\mathbf{z}_t; t)\|_2^2 \quad (25)$$

$$= \frac{1}{2} (\text{SNR}(s) - \text{SNR}(t)) \|\mathbf{z}_0 - \mathbf{z}_\theta(\mathbf{z}_t; t)\|_2^2, \quad (26)$$

where  $\text{SNR}(t) = \alpha_t^2 / \sigma_t^2$  and its physical meaning is signal-to-noise ratio.

After reparameterization of  $\mathbf{z}_t$ , the diffusion loss function becomes:

$$\mathcal{L}_{\text{VLB}}(\mathbf{z}_0) = \sum_{t=1}^T \mathbb{E}_{q(\mathbf{z}_t|\mathbf{z}_0)} [D_{KL}(q(\mathbf{z}_s|\mathbf{z}_t, \mathbf{z}_0) || p_\theta(\mathbf{z}_s|\mathbf{z}_t))] \quad (27)$$

$$= \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[ \sum_{t=1}^T (\text{SNR}(s) - \text{SNR}(t)) \|\mathbf{z}_0 - \mathbf{z}_\theta(\mathbf{z}_t; t)\|_2^2 \right]. \quad (28)$$

In practice, we follow Plaid to use the continuous-time diffusion formulation, where  $t \in [0, 1]$ , and we can express  $\mathcal{L}$  as a function of  $\tau$  with  $\tau \rightarrow 0$ :

$$\mathcal{L}_{\text{VLB}}(\mathbf{z}_0) = \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} \int_0^1 \left[ \frac{\text{SNR}(t - \tau) - \text{SNR}(t)}{\tau} \|\mathbf{z}_0 - \mathbf{z}_\theta(\mathbf{z}_t; t)\|_2^2 \right] dt, \quad (29)$$

and let  $\text{SNR}'(t)$  denote the derivative of the SNR function, this then gives:

$$\mathcal{L}_{\text{VLB}}(\mathbf{z}_0) = -\frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} \int_0^1 \text{SNR}'(t) \|\mathbf{z}_0 - \mathbf{z}_\theta(\mathbf{z}_t; t)\|_2^2 dt. \quad (30)$$

### A.3 Pre-trained SEDD

SEDD [37] is a discrete diffusion language model built based on discrete score matching [40], which generalizes score matching [45, 46] to the discrete data. We now denote  $x$  as a categorical random variable and the following derivation can be extended to a sequence of variable  $\mathbf{x}$  as well.

**Concrete score.** Instead of directly modeling  $p_\theta(x)$  to approximate original data distribution  $q(x)$ , the core idea of discrete score matching is to learn a quantity known as the *concrete score* [40] through a neural network:

$$s_\theta(x)_y = \frac{p_\theta(y)}{p_\theta(x)} = \frac{e^{f_\theta(y)}/Z}{e^{f_\theta(x)}/Z} = \frac{e^{f_\theta(y)}}{e^{f_\theta(x)}}, \quad (31)$$

which eliminates normalizing constant  $Z$  as in the energy-based model. In particular, this quantity is the categorical equivalent of the famous score function  $\nabla_x \log p$  in continuous space. Regarding the choice of  $y$ , if we model the ratio for every possible  $y$ , we would have  $V$  items given  $V$  as the dimension of  $x$ , and  $N^V$  items for  $\mathbf{x}$  given  $N$  as the sequence length of  $\mathbf{x}$ , which is computationally intractable. So we sparsify and only model "relevant" ratios based on whether  $y$  is "close" to  $x$ . These relevant positions will be denoted as  $y \sim x$ , e.g., all sentences  $y$  that differ from  $x$  with Hamming distance 1.

**Training objective.** Lou et al. [37] define a learning objective named *score entropy* to learn the neural network, which is given by:

$$\mathbb{E}_{x \sim q} \left[ \sum_{y \sim x} s_\theta(x)_y - \frac{q(y)}{q(x)} \log s_\theta(x)_y \right]. \quad (32)$$

Taking a derivative w.r.t.  $s$  and setting it to 0, we see that this occurs when  $s_\theta(x)_y = \frac{q(y)}{q(x)}$ , which can be easily checked to be globally optimal as the function is convex as a function of  $s$ . To handle the unknown term  $\frac{q(y)}{q(x)}$ , they further propose *denoising score entropy* motivated by *denoising score matching* [45] based on  $q(x_t) = \sum_{x_0} q_{t|0}(x_t|x_0)q_0(x_0)$ :

$$\mathbb{E}_{x_0 \sim q_0, t \sim U[0, T], x_t \sim q_{t|0}(x_t|x_0)} \left[ \sum_{y \sim x_t} s_\theta(x_t, t)_y - \frac{q_{t|0}(y|x_0)}{q_{t|0}(x_t|x_0)} \log s_\theta(x_t, t)_y \right], \quad (33)$$

where  $q_{t|0}(\cdot|x_0)$  is a perturbation of a base density  $q(\cdot)$  by a transition kernel, and the transition ratio  $\frac{q_{t|0}(y|x_0)}{q_{t|0}(x_t|x_0)}$  is known by design.

**Forward diffusion process.** The transition  $q_{t|0}(x_t|x_0)$  is a vector that represents a categorical distribution, and can be defined by a forward diffusion process  $q_{t|0}(x_t|x_0) = \exp(\bar{\sigma}(t)Q)_{x_0}$ , where  $\bar{\sigma}(t) \in \mathbb{R}_{\geq 0}$  is the cumulative noise  $\int_0^t \sigma(s)ds$  at timestep  $t$  with a value close to 0 when  $t$  is small and increasing when  $t$  growing. Lou et al. [37] use two standard transition matrices with special structures to implement matrix  $Q$  following prior work [1]:

$$Q^{\text{uniform}} = \begin{bmatrix} 1-V & 1 & \cdots & 1 \\ 1 & 1-V & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1-V \end{bmatrix} \quad (34)$$

$$Q^{\text{absorb}} = \begin{bmatrix} -1 & 0 & \cdots & 0 & 0 \\ 0 & -1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & -1 & 0 \\ 1 & 1 & \cdots & 1 & 0 \end{bmatrix} \quad (35)$$

One can view the above diffusion process by taking small  $\Delta t$  Euler steps and randomly sampling the resulting transitions:

$$q_{t+\Delta t|t}(x_{t+\Delta t} = y | x_t = x) \propto \delta_{xy} + Q_t(y, x)\Delta t + O(\Delta t), \quad (36)$$

where  $Q_t = \sigma(t)Q$ , and  $O(\Delta t)$  represents terms that tend to zero at a faster rate than  $\Delta t$ .

**Backward generative process.** To simulate the diffusion defined above, one can use the Euler strategy to derive the time reversal of the forward process:

$$q_{t-\Delta t|t}(x_{t-\Delta t} = y | x_t = x) \propto \delta_{xy} + R_t(y, x)\Delta t + O(\Delta t), \quad (37)$$

where  $R_t$  is the reverse transition rate matrix that can be derived using Bayes rule:  $R_t(y, x) = \frac{q_t(y)}{q_t(x)}Q_t(x, y)$ . Each column of  $R_t$  represents the transition probability from a token at timestep  $t$  to other tokens at timestep  $t - \Delta t$ . Let  $p_\theta(x_{t-\Delta t} = y | x_t = x) = q_{t-\Delta t|t}(x_{t-\Delta t} = y | x_t = x)$ , we have:

$$p_\theta(x_{t-\Delta t} = y | x_t = x) \propto \delta_{xy} + R_t^\theta(y, x)\Delta t + O(\Delta t), \quad (38)$$

where  $R_t^\theta(y, x) = \sum_{x_0} q_{0|t}(x_0|x) \frac{q_{t|0}(y|x_0)}{q_{t|0}(x|x_0)} Q_t(x, y) = s_\theta(x, t)_y Q_t(x, y)$ . For a sequence of random variables  $\mathbf{x}$ , this is inefficient because only one position is modified per step. A natural alternative has been to use  $\tau$ -leaping [12], which performs an Euler step at each position simultaneously.

#### A.4 Conditional ODE solver

The sampling of continuous diffusion models can be implemented by solving the diffusion ODEs [45, 46]. Specifically, sampling by diffusion ODEs needs to discretize the following ODE [46] with  $t$  changing from  $T$  to 0:

$$\frac{d\mathbf{z}_t}{dt} = f(t)\mathbf{z}_t + \frac{g^2(t)}{2\sigma_t^2} \epsilon_\theta(\mathbf{z}_t, t), \quad \mathbf{z}_T \sim \mathcal{N}(\mathbf{0}, \tilde{\sigma}^2 \mathbf{I}). \quad (39)$$

The *data prediction model*  $\mathbf{z}_\theta(\mathbf{z}_t, t)$  predicts the original data  $\mathbf{z}_0$  based on the noisy  $\mathbf{z}_t$ , and its relationship with  $\epsilon_\theta(\mathbf{z}_t, t)$  is given by  $\mathbf{z}_\theta(\mathbf{z}_t; t) := (\mathbf{z}_t - \sigma_t \epsilon_\theta(\mathbf{z}_t, t))/\alpha_t$  [28]. Therefore, the equivalent diffusion ODE w.r.t. the data prediction model  $\mathbf{z}_\theta$  is:

$$\frac{d\mathbf{z}_t}{dt} = \left( f(t) + \frac{g^2(t)}{2\sigma_t^2} \right) \mathbf{z}_t - \frac{\alpha_t g^2(t)}{2\sigma_t^2} \mathbf{z}_\theta(\mathbf{z}_t, t), \quad \mathbf{z}_T \sim \mathcal{N}(\mathbf{0}, \tilde{\sigma}^2 \mathbf{I}), \quad (40)$$

where the coefficients  $f(t) = \frac{d \log \alpha_t}{dt}$ ,  $g^2(t) = \frac{d\sigma_t^2}{dt} - 2 \frac{d \log \alpha_t}{dt} \sigma_t^2$  [28].

Given an initial value  $\mathbf{z}_s$  at time  $s > 0$  and denote  $\hat{\mathbf{z}}_\theta(\hat{\mathbf{z}}_\lambda, \lambda) := \mathbf{z}_\theta(\mathbf{z}_{t_\lambda(\lambda)}, t_\lambda(\lambda))$  as the change-of-variable form of  $\mathbf{z}_\theta$  for  $\lambda$ , the solution  $\mathbf{z}_t$  at time  $t \in [0, s]$  of diffusion ODEs in Eq. (40) is:

$$\mathbf{z}_t = \frac{\sigma_t}{\sigma_s} \mathbf{z}_s + \sigma_t \int_{\lambda_s}^{\lambda_t} e^{\lambda \hat{\mathbf{z}}_\theta(\hat{\mathbf{z}}_\lambda, \lambda)} d\lambda, \quad (41)$$

which can be proved by taking derivative w.r.t.  $t$  in Eq. (41):

$$\begin{aligned}\frac{d\mathbf{z}_t}{dt} &= \frac{d\sigma_t}{dt} \frac{\mathbf{z}_s}{\sigma_s} + \frac{d\sigma_t}{dt} \int_{\lambda_s}^{\lambda_t} e^{\lambda} \hat{\mathbf{z}}_{\theta}(\hat{\mathbf{z}}_{\lambda}, \lambda) d\lambda + \frac{d\lambda_t}{dt} \sigma_t e^{\lambda_t} \hat{\mathbf{z}}_{\theta}(\hat{\mathbf{z}}_{\lambda_t}, \lambda_t) \\ &= \frac{d\sigma_t}{dt} \frac{\mathbf{z}_t}{\sigma_t} + \frac{d\lambda_t}{dt} \sigma_t e^{\lambda_t} \hat{\mathbf{z}}_{\theta}(\hat{\mathbf{z}}_{\lambda_t}, \lambda_t) \\ &= \left( f(t) + \frac{g^2(t)}{2\sigma_t^2} \right) \frac{\mathbf{z}_t}{\sigma_t} - \frac{\alpha_t g^2(t)}{2\sigma_t^2} \mathbf{z}_{\theta}(\mathbf{z}_t, t),\end{aligned}$$

and this gives us the exact formulation as in Eq. (41).

Based on Eq. (41), the aim of an ODE solver is to approximate the exact solution at time  $t_i$  given the previous value  $\mathbf{z}_{t_{i-1}}$  at time  $t_{i-1}$ . Denote  $\mathbf{z}_{\theta}^{(n)}(\lambda) := \frac{d^n \hat{\mathbf{z}}_{\theta}(\mathbf{z}_{\lambda}, \lambda)}{d\lambda^n}$  as the  $n$ -th order total derivatives of  $\mathbf{z}_{\theta}$  w.r.t. logSNR  $\lambda$ . Lu et al. [38, 39] show that by taking the  $(k-1)$ -th Taylor expansion ( $k \geq 1$ ) at  $\lambda_{t_{i-1}}$  for  $\mathbf{z}_{\theta}$  w.r.t.  $\lambda \in [\lambda_{t_{i-1}}, \lambda_{t_i}]$  and substitute it into Eq. (41) with  $s = t_{i-1}$  and  $t = t_i$ , we have

$$\mathbf{z}_{t_i} = \frac{\sigma_{t_i}}{\sigma_{t_{i-1}}} \mathbf{z}_{t_{i-1}} + \underbrace{\sigma_{t_i} \sum_{n=0}^{k-1} \mathbf{z}_{\theta}^{(n)}(\hat{\mathbf{z}}_{\lambda_{t_{i-1}}}, \lambda_{t_{i-1}})}_{\text{estimated}} \underbrace{\int_{\lambda_{t_{i-1}}}^{\lambda_{t_i}} e^{\lambda} \frac{(\lambda - \lambda_{t_{i-1}})^n}{n!} d\lambda}_{\text{analytically computed}} + \underbrace{\mathcal{O}(h_i^{k+1})}_{\text{omitted}}, \quad (42)$$

where the integral  $\int e^{\lambda} \frac{(\lambda - \lambda_{t_{i-1}})^n}{n!} d\lambda$  can be analytically computed by integral-by-parts. Therefore, to design a  $k$ -th order ODE solver, we only need to estimate the  $n$ -th order derivatives  $\mathbf{z}_{\theta}^{(n)}(\lambda_{t_{i-1}})$  for  $n \leq k-1$  after omitting the  $\mathcal{O}(h_i^{k+1})$  high-order error terms. For  $k=1$ , Eq. (42) becomes (after omitting the  $\mathcal{O}(h_i^{k+1})$  terms)

$$\mathbf{z}_{t_i} = \frac{\sigma_{t_i}}{\sigma_{t_{i-1}}} \mathbf{z}_{t_{i-1}} + \sigma_{t_i} \mathbf{z}_{\theta}(\mathbf{z}_{t_{i-1}}, t_{i-1}) \int_{\lambda_{t_{i-1}}}^{\lambda_{t_i}} e^{\lambda} d\lambda = \frac{\sigma_{t_i}}{\sigma_{t_{i-1}}} \mathbf{z}_{t_{i-1}} - \alpha_{t_i} (e^{-h_i} - 1) \mathbf{z}_{\theta}(\mathbf{z}_{t_{i-1}}, t_{i-1}), \quad (43)$$

where  $h_i := \lambda_{t_i} - \lambda_{t_{i-1}}$  for  $i = 1, \dots, T$ .

Since DoT is conditionally trained with partial nosing, we introduce a conditional form of Eq. (43) when adapting the above ODE solver into the inference stage. For  $k=1$ , this is given by:

$$\mathbf{y}_{t_i} = \frac{\sigma_{t_i}}{\sigma_{t_{i-1}}} \mathbf{y}_{t_{i-1}} - \alpha_{t_i} (e^{-h_i} - 1) \tilde{\mathbf{z}}_{\theta}(\mathbf{z}_{t_{i-1}}, t_{i-1}),$$

where  $\mathbf{z}_{t_{i-1}} = [\mathbf{x}; \mathbf{y}_{t_{i-1}}]$  and  $\tilde{\mathbf{z}}_{\theta}(\mathbf{z}_t, t)$  is used to denote the fractions of recovered  $\mathbf{z}_0$  corresponding to  $\mathbf{y}_0$ .

## B Experiment Details

### B.1 Dataset Statistics

We list the statistics of our used datasets in Table 4. For the digit multiplication datasets and GSM8K dataset, we use processed datasets from Implicit CoT<sup>6</sup> [7]. For boolean logic task, we construct the training and test dataset using the method from DyVal<sup>7</sup> [68]. All datasets contain 1000 test examples except GSM8K, which contains 1319 examples.

### B.2 Details of Baselines

When fine-tuning GPT2, we train 40 epochs using the learning rate of  $1e-4$  for boolean logic and  $5e-4$  for others. During inference, we use greedy decoding for single decoding. For self-consistency, following the original paper [52], we apply temperature sampling with  $T = 0.5$  and truncated at the top- $k$  ( $k = 40$ ) tokens with the highest probability for diverse generation. All GPT2-based models

<sup>6</sup>[https://github.com/da03/implicit\\_chain\\_of\\_thought](https://github.com/da03/implicit_chain_of_thought)

<sup>7</sup><https://github.com/microsoft/promptbench/blob/main/examples/dyval.ipynb>

Table 4: Training set size, average number of tokens in the input, intermediate, and output texts respectively when using Plaid tokenizer on the validation set and average number of rationales.

Dataset	Size	#Input token	#Intermediate token	#Output token	#Rationales
4x4	808k	16	84	15	4
5x5	808k	20	137	19	5
Boolean logic	99k	112	134	3	10
GSM8K-Aug	378k	61	34	2	2.7

use GPT2Tokenizer with vocabulary size of 50257. All datasets are trained using sequence length of 256 except boolean logic, which uses 384 length.

Note in Table 1, we compare Plaid DoT with the fine-tuned GPT2 small, given that the Plaid 1B [18] model exhibits similar perplexity to GPT2 small. This might put our Plaid DoT model at a disadvantage in terms of inference speed, as the parameters of Plaid 1B are nearly  $10\times$  greater than those of GPT2 small.

For Transformer-scratch baseline [51], we use 6 transformer encoder layers and 6 transformer decoder layers. We employ the tokenizer from bert-base-uncased with a vocabulary size of 30522. The learning rate is set to  $1e-5$ , and we train for 60k steps with a batch size of 128.

For ChatGPT, we use OpenAI api<sup>8</sup> with the following prompt in 5-shot.

<p>Answer the final question following the format of the given examples.</p> <p>Example problems:</p> <p>Q: {query} A: {answer} ... Question to answer: Q:</p>
--

Table 5: Prompt for ChatGPT.

Please note that the throughput of ChatGPT in Table 1 only measures the response speed of ChatGPT and does not represent the actual generation speed of the model. As a blackbox commercial product, ChatGPT may employ various optimization techniques to speedup generating responses to enhance user experiences.

### B.3 DoT Implementation Details

We conduct all the experiments on NVIDIA V100-32G GPUs, and we use 8 GPUs for training and sampling. We resort to half precision (fp16) instead of bfloat16 (bf16) as V100 GPU doesn't support bf16, and we don't observe any number explosion. By default, we train DoT from scratch on three datasets respectively, including the four-digit ( $4 \times 4$ ), five-digit ( $5 \times 5$ ) multiplication datasets, and the GSM8k dataset. Additionally, we fine-tune the pre-trained model Plaid-1B on the GSM8K dataset with DoT to explore its effectiveness further.

For DoT trained from scratch. We use 12 layers of transformer and bert-base-uncased vocabulary. We preprocess the four-digit ( $4 \times 4$ ) and five-digit ( $5 \times 5$ ) multiplication datasets to prepare for the training process of the DoT multi-path variant, and sampling from it. The learning rate is  $1e-4$  and we train for 60k steps with the batch size of 128 and max sequence length of 128. For digit

<sup>8</sup><https://platform.openai.com/docs/api-reference>

multiplication in Table 1, we use sampling step  $T = 1$  to achieve high throughput while keeping the accuracy. For boolean logic dataset, we use  $T = 2$ .

For DoT fine-tuned from Plaid, we set the training steps of the DoT and multi-pass DoT to be 120k and 30k respectively, as we find more training steps will lead to performance degradation. The learning rate is set to  $1\text{e-}4$  for boolean logic and  $3\text{e-}4$  for other datasets. The max sequence length is set to 384 for the boolean logic dataset and 256 for others. We use Adam optimizer [29]. During tokenization, we use Plaid’s tokenizer and we treat all the digits as individual tokens. During training, we set  $\epsilon_{min}$  to be 0.95 as we find decreasing the probability of oracle demonstration hinders model training. We choose glancing sampling  $\gamma = 0.01$  and self consistency  $m = 20$ . Following Gulrajani and Hashimoto [18], we also adopt self-conditioning [5] during training. During inference, we set the scoring temperature to 0.5 to sharpen the predicted noise distribution. We also use soft logits with a temperature of 0.5 to produce more diverse samples. By default, we use sampling step  $T = 64$  to ensure accuracy. Training DoT and DoT require 29h and 10h, respectively. For DoT trained from SEDD, we set the training steps of the DoT and multi-pass DoT to be 200k, with other parameters being the same as when training Plaid. For all the experiments, we have verified the statistical significance by running them multiple times.

#### B.4 Additional Results

Table 6: Comparison to larger AR models.

	Params	Accuracy
GPT-2-medium CoT	355M	43.9
Mistral CoT	7B	68.8
Llama CoT	7B	59.0
SEDD-medium DoT <sup>MP</sup>	424M	53.5

**Comparison to larger open language models.** We compare our model with LoRA fine-tuning of AR LLMs on the same GSM-Aug dataset, which is listed in Table 6. Please note that the current diffusion pretrained model is much smaller than Llama 7B, so this comparison is not fair and we just list them for reference. We have validated that our DoT is better than the same scale autoregressive model GPT-2, which shares a similar architecture with Llama. We believe that further exploration of diffusion language models will lead to larger models that can compete with current LLMs, allowing DoT to achieve results more comparable to Llama.

Table 7: Comparison between DoT and no-DoT (Answer-only).

	Accuracy
GPT-2-small Answer-only	13.3
GPT-2-small CoT	39.0
Plaid Answer-only	12.4
Plaid DoT <sup>MP</sup>	37.7
SEDD-small Answer-only	29.1
SEDD-small DoT <sup>MP</sup>	43.2

**Comparison with no-DoT finetune.** We conduct the answer-only setting to further validate the effectiveness of DoT. The results in Table 7 reveal that fine-tuning diffusion models solely with answer data leads to inferior performance compared to DoT, mirroring the degradation of AR models in the absence of CoT.

**Throughput Comparison.** We have shown how  $T$  affects performance on grade school math in Figure 3, and here we also show how  $T$  affects throughput for Plaid DoT<sup>MP</sup>, as in Table 8. The relationship between throughput and  $T$  appears to be nearly linear.

**Comparison of the reasoning paths between DoT and DoT<sup>MP</sup>.** We observe that DoT<sup>MP</sup> outperforms DoT in correctness regarding the reasoning paths, while DoT slightly excels in diversity

Table 8: Throughput comparison when increasing the number of timesteps  $T$  for Plaid DoT<sup>MP</sup>.

T	Accuracy	Throughput
1	18.2	6.6
2	35.9	3.4
4	36.7	1.7
8	36.4	0.9
16	36.1	0.4
32	37.4	0.2
64	37.7	0.1
128	37.7	.05

as depicted in Figure 4(b). Below we show some examples where DoT<sup>MP</sup> can predict the correct reasoning path while DoT fails:

**Query:** The Kennel house keeps 3 German Shepherds and 2 Bulldogs. If a German Shepherd consumes 5 kilograms of dog food and a bulldog consumes 3 kilograms of dog food per day. How many kilograms of dog food will they need in a week?

**DoT:** «3\*5=15» «7\*3=21» «15+21=36» ##### 36

**DoT<sup>MP</sup>:** «3\*5=15» «2\*3=6» «15+6=21» «21\*7=147» ##### 147

**Query:** Skyler has 100 hats on his hand with the colors red, blue, and white. Half of the hats are red, 3/5 of the remaining hats are blue, and the rest are white. How many white hats does Skyler have?

**DoT:** «1/2\*100=50» «3/5\*50=30» «100-30=70» ##### 70

**DoT<sup>MP</sup>:** «100/2=50» «100-50=50» «50\*3/5=30» «50-30=20» ##### 20

## B.5 Other Attempts

For the ablation design for DoT fine-tuning in Table 2, we have tried to fine-tune a decoder-only autoregressive language model (i.e., GPT2 here), where we only change the base model from Plaid 1B to GPT2 large, remove the causal mask and keep all other diffusion training settings the same with the Plaid fine-tuning. In this setting, even though the model is formulated and trained in the diffusion manner, it still can not predict the right format of answers. This experiment may indicate that a pre-trained diffusion model is necessary for the further fine-tuning of downstream tasks.

Regarding datasets, we also try to mix up four-digit ( $4 \times 4$ ) and five-digit ( $5 \times 5$ ) multiplication datasets for training and testing, considering that the number of rationales is different in these two tasks. As for the result, the trained model learns when to conclude the computation and can attain 100% accuracy.

## C Discussion about base models

Our DoT approach is constrained by the pre-training and fine-tuning paradigm due to the not-strong-enough base models. This lags behind the current trend of instruction-tuning LLMs and pursuing the generalization of LMs across various tasks. Nevertheless, considering the pre-trained diffusion models are still in their early stages and the lack of scaled pre-trained diffusion models, our study is a preliminary exploration to show the potential of diffusion models for reasoning tasks, and we believe that with more powerful pre-trained diffusion models and post-instruction tuning, DoT can attain the generalization capabilities of today’s LLMs and yield further advantages.

## D Boarder Impacts

Our work contributes to the understanding of denoising generative models and enhances their generation capabilities within certain discrete text reasoning datasets. The proposed DoT with

diffusion language models challenges autoregressive models with CoT, achieving competitive performance. While there is still a large gap with modern large autoregressive language models such as ChatGPT, we believe DoT can benefit more with future work on scaling diffusion language models. However, we acknowledge that deep generative models, as powerful tools for learning from unstructured data, can have detrimental societal impacts if misused. Specifically, these models can facilitate the spread of misinformation by reducing the resources required to create realistic fake content. Additionally, the generated samples from these models accurately reflect the statistics of their training datasets. Consequently, if these samples are interpreted as objective truth without considering the inherent biases present in the original data, they can perpetuate discrimination against minority groups.



## NeurIPS Paper Checklist

### 1. Claims

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

Answer: [Yes]

Justification: Please see abstract and introduction.

Guidelines:

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

### 2. Limitations

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

Answer: [Yes]

Justification: We discuss the limitations in section 6.

Guidelines:

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

### 3. Theory Assumptions and Proofs

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

Answer: [NA]

Justification: This paper does not include theoretical results.

Guidelines:

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

#### 4. Experimental Result Reproducibility

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

Answer: [\[Yes\]](#)

Justification: We provide the details of the dataset, training infrastructure, and implementations in Appendix B.

Guidelines:

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

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [\[Yes\]](#)

Justification: We have uploaded the code to reproduce our results.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

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

Answer: [\[Yes\]](#)

Justification: We provide the implementation details in Section 4.1 and Appendix B. We also provide the official code implementation in ...

Guidelines:

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

## 7. Experiment Statistical Significance

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

Answer: [\[Yes\]](#)

Justification: Please see Appendix B.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

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

Answer: [\[Yes\]](#)

Justification: Please see Appendix B.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code Of Ethics

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

Answer: [\[Yes\]](#)

Justification: We have reviewed the NeurIPS Code of Ethics.

Guidelines:

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

## 10. Broader Impacts

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

Answer: [\[Yes\]](#)

Justification: Please see Appendix D.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

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

Answer: [NA]

Justification: The paper poses no such risks

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

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

Answer: [Yes]

Justification: Please see Section 4.1 and Appendix B.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 13. **New Assets**

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

Answer: [\[Yes\]](#)

Justification: The license detail is provided in the uploaded code.

Guidelines:

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

### 14. **Crowdsourcing and Research with Human Subjects**

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

Answer: [\[NA\]](#)

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

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

### 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

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

Answer: [\[NA\]](#)

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

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