

812 A Transformer and Remasking Step

Algorithm 2 Diffusion Step

Require: Transformer \mathcal{N}_n , full output length n , prompt \mathbf{p} , input prompt length m , block length d , current diffusion step i , total diffusion steps T , vocabulary V .

```

1:  $\mathbf{x} \leftarrow \mathbf{p} \cdot \perp^d$  ▷ Pad the input prompt where  $n = m + d$ 
2:  $\mathbf{v}_1 \dots \mathbf{v}_n \leftarrow \mathcal{N}_n(\mathbf{x})$  ▷  $\mathbf{v}_i \in \mathbb{R}_+^{|V|}$  output distribution at position  $i$ 
3:  $\mathbf{l} \leftarrow \text{RemaskPositions}(\mathbf{v}_{m+1}, \dots, \mathbf{v}_{m+d}, i, T)$  ▷ Decides which positions to remask
4: for  $j \in \mathbf{l}$  do
5:    $\mathbf{v}_j \leftarrow \mathbf{0}$ 
6:    $\mathbf{v}_j[\perp] \leftarrow 1$  ▷ Set probability of all tokens except  $\perp$  to 0.
7:  $\mathbf{r} \leftarrow D_{m,n}(\mathbf{v}_1 \dots \mathbf{v}_n)$  ▷ Decoding that outputs response with first  $m$  tokens are input prompt  $\mathbf{p}$ 
8: return  $\mathbf{r}$ 

```

813 We describe the two key components of a single diffusion step: a) **Transformer step:** Computes
 814 the output distribution over all tokens in the vocabulary (line 2 Algo. 2). b) **Remasking step:** Based
 815 on the output from the transformer step, it greedily decides which token positions to mask. The
 816 remasking step can be viewed as updating the output distribution such that, at the masked positions,
 817 the mask token \perp is assigned probability 1, while all other tokens receive probability 0 (lines 3 – 6
 818 Algo. 2). Popular greedy remasking strategies include (line 3 Algo. 2): (i) *Random*: Masks tokens at
 819 randomly selected positions Nie et al. [2025]. (ii) *Top token probability*: Masks positions where the
 820 top-predicted token has the lowest probability Nie et al. [2025]. (iii) *Entropy-based*: Computes the
 821 entropy of the output distribution at each position and masks the positions with the highest entropy
 822 Ye et al. [2025].

823 The number of token positions to remask at the i -th step typically depends on the total number of
 824 diffusion steps T and the block length d . At step 0, all d positions are masked, and the number of
 825 masked tokens decreases linearly to 0 over T steps. Thus, at the i -th step, the number of masked
 826 tokens is given by $\left\lfloor \frac{d \times (T-i)}{T} \right\rfloor$.

827 B Proofs

828 **Proposition 4.1.** [Correctness] Given any regular expression \mathcal{R} , input prompt $\mathbf{p} \in V^m$, block length
 829 d , output distribution $\mathcal{D}_{m+d} = \mathbf{v}_1 \dots \mathbf{v}_{m+d}$, if $L_P(\mathcal{R}) \cap (V \setminus \perp)^d \neq \{\}$ and $\mathbf{r} \sim \mathbf{v}_{m+1} \dots \mathbf{v}_{m+d}$ be
 830 the decoded string, then $\exists \mathbf{x} \in V^* . (\mathbf{x} \in \mathcal{S}(\mathbf{r})) \wedge (\mathbf{x} \in L_P(\mathcal{R}))$ holds.

831 *Proof.* We assume that $\exists \mathbf{x} \in L_P(\mathcal{R}) \cap (V \setminus \perp)^d \wedge (P(\mathbf{x} | \mathbf{v}_{m+1} \dots \mathbf{v}_{m+d}) \geq 0)$ then the decoded
 832 string \mathbf{r} satisfy the soundness property (see Definition 3.2). In other words, if there is at least one
 833 fully unmasked valid prefix with non-zero probability then DINGO retrieves a valid string.

834 We show this by induction on the position of tokens. Before moving to the proof, we first define
 835 extended transition function δ^* when $\delta : Q \times V \rightarrow 2^Q$ outputs a set of states instead of single
 836 state due to mask token \perp . In this case, for any string $\mathbf{w} \in V^*$, $\delta^*(\mathbf{w}, q_0)$ represents the state of
 837 reachable states starting from q_0 . This can be defined as $\delta^*(\{\}, q_0) = \{q_0\}$ and $\delta^*(t_1 \dots t_{m+1}, q_0) =$
 838 $\cup_{q \in \delta^*(t_1 \dots t_m, q_0)} \delta(q, t_{m+1})$.

839 1. Let $0 \leq i \leq d$, and let $t_1 \dots t_i \in V^i$ denote any token sequence with positive probability
 840 mass $\prod_{j=1}^i \mathbf{v}_{m+j}[t_j] > 0$. Let $q \in \delta^*(t_1 \dots t_i, q_0)$. Then, $W[i, q] > 0$. We prove this using
 841 induction on i .

842 (a) Base case $i = 0$: For empty strings only start state q_0 is reachable. DINGO initializes
 843 $W[0, q_0] = 1 > 0$ and for all $q \neq q_0$, $W[0, q] = 0$. (lines 1 – 3 in Algo. 1).

844 (b) Inductive Step: At position $i + 1$, let $t_1 \dots t_{i+1} \in V^{i+1}$ s.t. $\prod_{j=1}^{i+1} \mathbf{v}_{m+j}[t_j] > 0$. Let
 845 $q' \in \delta^*(t_1 \dots t_i, q_0)$ and $q \in \delta(q', t_{i+1})$. By the inductive hypothesis, for all such q'

846

$W[i, q'] > 0$. Recall,

$$V_{i+1}(q, q') = \begin{cases} \max_{t \in V} \mathbf{v}_{m+i+1}(t) \text{ s.t. } q \in \delta(q', t) \\ 0 \text{ if } q, q' \text{ are not connected} \end{cases} \quad W[i+1, q] = \max_{q' \in Q} W[i, q'] \times V_{i+1}(q, q')$$

847

Thus, $V_{i+1}(q, q') \geq \mathbf{v}_{m+i+1}(t_{i+1}) > 0$ which implies $W[i, q'] \times V_{i+1}(q, q') > 0$.

848

Therefore, $W[i+1, q] = \max_{q' \in Q} W[i, q'] \times V_{i+1}(q, q') > 0$.

849

2. Since $L_P(\mathcal{R}) \cap (V \setminus \perp)^d \neq \{\}$ by assumption, there exists some $\mathbf{y} \in L_P(\mathcal{R}) \cap (V \setminus \perp)^d$. By the Definition 2.6, $q_l = \delta_t^*(\mathbf{y}, q_0) \in Q_l$. From the induction above, $W[d, q_l] > 0$. From line 16 in Algo. 1, $q_{max} = \arg \max_{q \in Q_l} W[d, q]$. Thus, by the definition of $\arg \max$, $W[d, q_{max}] \geq W[d, q_l] > 0$.

850

851

852

853

854

855

856

857

3. In lines 20-22 in Algo. 1), DINGO reconstructs a d -length sequence $r = t_1 \dots t_d \in V^d$ such that $q_{max} \in \delta^*(r, q_0)$. For any $t_j \in r$, if $t_j = \perp$, choose any token $\tau_j \in (V \setminus \perp)$ satisfying $\delta_t(q_{j-1}, \tau_j) = q_j$ where $q_j = \delta_t^*(t_1 \dots t_j, q_0)$. By definition of δ_\perp , τ_j exists. Substituting every \perp in this manner yields, by Definition 3.1, $\mathbf{x} = \mathbf{x}_1 \dots \mathbf{x}_d \in (V \setminus \perp)^d$. $\mathbf{x} \in \mathcal{S}(\mathbf{r})$. $\delta_t^*(\mathbf{x}, q_0) = q_{max}$. From above, $W[d, q_{max}] > 0$.

858

859

4. Since $q_{max} \in Q_l$, by Definition 2.6, $\exists w \in \Sigma^*$ s.t. $\delta^*(w, q_{max}) \in F$. Equivalently, $\mathbf{x} \cdot w \in L(\mathcal{R})$, hence $\mathbf{x} \in L_P(\mathcal{R})$.

860

□

861

862

863

864

Proposition 4.2. [Optimality] Given any regular expression \mathcal{R} , input prompt $\mathbf{p} \in V^m$, block length d , output distribution $\mathcal{D}_{m+d} = \mathbf{v}_1 \dots \mathbf{v}_{m+d}$, if $L_P(\mathcal{R}) \cap (V \setminus \perp)^d \neq \{\}$ and $\mathbf{r}^* \sim \mathbf{v}_{m+1} \dots \mathbf{v}_{m+d}$ be the decoded string, then for any valid string \mathbf{r}' satisfying $\exists \mathbf{x} \in V^* . (\mathbf{x} \in \mathcal{S}(\mathbf{r}')) \wedge (\mathbf{x} \in L_P(\mathcal{R}))$, $P(\mathbf{r}' \mid \mathbf{v}_{m+1} \dots \mathbf{v}_n) \leq P(\mathbf{r}^* \mid \mathbf{v}_{m+1} \dots \mathbf{v}_n)$.

865

866

867

868

Proof. 1. First, we show that $P(\mathbf{r}^* \mid \mathbf{v}_{m+1} \dots \mathbf{v}_n) = W[d, q_{max}]$, or equivalently $\prod_{j=1}^d \mathbf{v}_{m+j}[\mathbf{r}_j^*] = W[d, q_{max}]$. Let $\mathbf{r}^* = \mathbf{r}_1^* \dots \mathbf{r}_d^*$ and $0 \leq i \leq d$. We prove by induction on i that if DINGO's backtracking (lines 19 – 23 in Algo 1) has brought us to state $q \in Q$ at position i , then $W[i, q] = \prod_{j=1}^i \mathbf{v}_{m+j}[\mathbf{r}_j^*]$.

869

870

871

872

873

- (a) Base case $i = 0$: $W[0, q_0] = 1 = \prod_{j=1}^0 \mathbf{v}_{m+j}[\mathbf{r}_j^*]$.
 (b) Inductive Step: At position i , let $q', \mathbf{r}_i^* = Pr[i, q]$ (line 21 in Algo 1). From lines 14 – 15 in Algo 1, $W[i, q] = W[i-1, q'] \times \mathbf{v}_{m+i}(\mathbf{r}_i^*)$. By the inductive hypothesis, $W[i-1, q'] = \prod_{j=1}^{i-1} \mathbf{v}_{m+j}[\mathbf{r}_j^*]$. Thus, $W[i, q] = \prod_{j=1}^{i-1} \mathbf{v}_{m+j}[\mathbf{r}_j^*] \times \mathbf{v}_{m+i}(\mathbf{r}_i^*) = \prod_{j=1}^i \mathbf{v}_{m+j}[\mathbf{r}_j^*]$.

874

875

Let $q_d \in \delta^*(\mathbf{r}_1^* \dots \mathbf{r}_d^*, q_0)$. Since $q_d = q_{max}$ (line 19 in Algo 1), $W[d, q_{max}] = \prod_{j=1}^d \mathbf{v}_{m+j}[\mathbf{r}_j^*] = P(\mathbf{r}^* \mid \mathbf{v}_{m+1} \dots \mathbf{v}_n)$.

876

877

878

2. We show that for every valid string $\mathbf{r}' = \mathbf{r}_1' \dots \mathbf{r}_d'$ satisfying $\exists \mathbf{x} \in V^* . (\mathbf{x} \in \mathcal{S}(\mathbf{r}')) \wedge (\mathbf{x} \in L_P(\mathcal{R}))$, $\prod_{j=1}^d \mathbf{v}_{m+j}[\mathbf{r}_j'] \leq W[d, q_{max}]$. Let $0 \leq i \leq d$ and $q \in \delta^*(\mathbf{r}_1' \dots \mathbf{r}_i', q_0)$. We show that $\prod_{j=1}^i \mathbf{v}_{m+j}[\mathbf{r}_j'] \leq W[i, q]$ using induction on i .

879

880

881

- (a) Base case $i = 0$: $W[0, q_0] = 1 = \prod_{j=1}^0 \mathbf{v}_{m+j}[\mathbf{r}_j']$.
 (b) Inductive Step: At position $i+1$, let $q' \in \delta^*(\mathbf{r}_1' \dots \mathbf{r}_i', q_0)$ and $q \in \delta(q', \mathbf{r}_{i+1}')$. By the inductive hypothesis, $\prod_{j=1}^i \mathbf{v}_{m+j}[\mathbf{r}_j'] \leq W[i, q']$. Recall,

$$V_{i+1}(q, q') = \begin{cases} \max_{t \in V} \mathbf{v}_{m+i+1}(t) \text{ s.t. } q \in \delta(q', t) \\ 0 \text{ if } q, q' \text{ are not connected} \end{cases} \quad W[i+1, q] = \max_{q' \in Q} W[i, q'] \times V_{i+1}(q, q')$$

882

883

Thus, $\mathbf{v}_{m+i+1}(\mathbf{r}_{i+1}') \leq V_{i+1}(q, q')$. Hence, $\prod_{j=1}^{i+1} \mathbf{v}_{m+j}[\mathbf{r}_j'] = \prod_{j=1}^i \mathbf{v}_{m+j}[\mathbf{r}_j'] \times \mathbf{v}_{m+i+1}(\mathbf{r}_{i+1}') \leq W[i, q'] \times V_{i+1}(q, q') \leq W[i+1, q]$.

884 Let $q_d \in \delta^*(\mathbf{r}_1' \dots \mathbf{r}_d', q_0)$. Since $\mathbf{x} \in V^* \cdot (\mathbf{x} \in \mathcal{S}(\mathbf{r}')) \wedge (\mathbf{x} \in L_P(\mathcal{R}))$, $q_d \in Q_l$. From
885 line 16 in Algo. 1, $q_{max} = \arg \max_{q \in Q_l} W[d, q]$. Thus, by the definition of $\arg \max$,
886 $W[d, q_d] \leq W[d, q_{max}]$. From the inductive hypothesis above, $\prod_{j=1}^d \mathbf{v}_{m+j}[\mathbf{r}_j'] \leq$
887 $W[d, q_d] \leq W[d, q_{max}]$.

888 3. Hence, $P(\mathbf{r}' \mid \mathbf{v}_{m+1} \dots \mathbf{v}_n) = \prod_{j=1}^d \mathbf{v}_{m+j}[\mathbf{r}_j'] \leq W[d, q_{max}] = \prod_{j=1}^d \mathbf{v}_{m+j}[\mathbf{r}_j^*] =$
889 $P(\mathbf{r}^* \mid \mathbf{v}_{m+1} \dots \mathbf{v}_n)$.

890 □

891 C Time complexity analysis of parallelized DINGO DP

Algorithm 3 DINGO DP

Require: q_0 , block length d , probability vectors $\mathbf{v}_1, \dots, \mathbf{v}_d$ for the current block, Q_l, Q, δ .

- 1: $W[0, q] \leftarrow 0$ for all $(q \in Q) \wedge (q \neq q_0)$
- 2: $W[0, q_0] \leftarrow 1$
- 3: $Pr[0, q] \leftarrow (\text{None}, \text{None})$ for all $(q \in Q)$ ▷ Initialization of the DP
- 4: $V_i \leftarrow \{\}$ for all $i \in \{1, \dots, d\}$ ▷ maximum token probability transtion ($q' \rightarrow q$) at position i
- 5: $T_i \leftarrow \{\}$ for all $i \in \{1, \dots, d\}$ ▷ token for the maximum probability transition ($q' \rightarrow q$)
- 6: **for** $i \in \{1, \dots, d\}$ **do** ▷ The computation along all d can be parallelized
- 7: # Parallelize for each $\{1, \dots, d\}$
- 8: **for** $(q \in Q)$ **do**
- 9: **for** $t \in V$ **do**
- 10: $q' \leftarrow \delta(q, t)$
- 11: $V_i(q, q'), T_i(q, q') \leftarrow \text{MaxTransition}(\mathbf{v}_i, t, q, q')$
- 12: **for** $i \in \{1, \dots, d\}$ **do** ▷ DP computation loop
- 13: **for** $(q \in Q) \wedge (q' \in Q)$ **do**
- 14: **if** $W[i, q] < W[i-1, q'] \times V_i(q, q')$ **then**
- 15: $W[i, q] \leftarrow W[i-1, q'] \times V_i(q, q')$ ▷ Update maximum probability path to q
- 16: $Pr[i, q] \leftarrow (q', T_i(q, q'))$ ▷ Update the parents accordingly
- 17: $q_{max} \leftarrow \arg \max_{q \in Q_l} W[d, q]$
- 18: **if** $W[d, q_{max}] = 0$ **then** ▷ No valid prefixes
- 19: **return** None, q_{max}
- 20: $\mathbf{r}^* \leftarrow \{\}, q_{curr} \leftarrow q_{max}$
- 21: **for** $i \in \{d, \dots, 1\}$ **do** ▷ Decoding the optimal string \mathbf{r}^*
- 22: $q_{curr}, t \leftarrow Pr[i, q_{curr}]$
- 23: $\mathbf{r}^* \leftarrow \mathbf{r}^* \cdot t$
- 24: **return** $\text{reverse}(\mathbf{r}^*), q_{max}$

892 The parallelism step at line 6 in Algo. 3 can be efficiently implemented using popular frameworks
893 like PyTorch. With parallelism, the computational depth (i.e., the minimum number of sequential
894 steps) reduces to $O(\max(|Q|^2, |Q| \times |V|) + |Q|^2 \times d)$. For regular expressions, where the number
895 of states $|Q|$ is a small constant, the computational depth becomes $O(|V| + d)$, which is linear in
896 both the vocabulary size $|V|$ and the block length d .

897 D Semi-Autoregressive

898 In the semi-autoregressive setup, given an input $\mathbf{p} \in V^m$, the output $\mathbf{o} \in V^{m+d \times k}$ is generated over
899 k blocks, where each block is computed via a call to the single block diffusion model. The output of
900 the i -th diffusion model call is $\mathbf{x}_i = \mathcal{L}_{m_i, n_i}(\mathbf{x}_{i-1})$, with $\mathbf{x}_0 = \mathbf{p}$ and the final output $\mathbf{o} = \mathbf{x}_k$. The
901 input and output lengths for each block are defined as $m_i = m + (i-1) \times d$ and $n_i = m + i \times d$ for
902 all $1 \leq i \leq k$.

Algorithm 4 Semi-Autoregressive diffusion LLM Generation

Require: diffusion LLM \mathcal{L} , prompt \mathbf{p} , answer length n , block length d , diffusion steps T , vocabulary V , number of blocks k .

```
1:  $\mathbf{x} \leftarrow \mathbf{p}$  ▷ Initialize  $\mathbf{x}$  with input prompt  $\mathbf{p}$ 
2:  $\mathbf{r} \leftarrow \{\}$  ▷ Initialize the output string
3: for  $i \in \{1, \dots, k\}$  do
4:    $\mathbf{x} \cdot \mathbf{r}_i \leftarrow \text{Diffusion}(\mathbf{x}, m + (i - 1) \times d, d, T, V)$  ▷  $\mathbf{r}_i \in V^d$  is  $i$ -th output block
5:    $\mathbf{r} \leftarrow \mathbf{r} \cdot \mathbf{r}_i$ 
6:    $\mathbf{x} \leftarrow \mathbf{x} \cdot \mathbf{r}_i$  ▷ Compute the input prompt for the next block
7: Return  $\mathbf{r}$ 
```

Algorithm 5 Semi-Autoregressive Constrained diffusion LLM Generation

Require: diffusion LLM \mathcal{L} , prompt \mathbf{p} , answer length n , block length d , diffusion steps T , vocabulary V , number of blocks k , regular expression \mathcal{R} .

```
1:  $q_0, Q_l, \delta \leftarrow \text{PreProcess}(\mathcal{R})$  ▷ Pre-compute the dfa start state, live states and  $\delta$ 
2:  $\mathbf{x} \leftarrow \mathbf{p}$  ▷ Initialize  $\mathbf{x}$  with input prompt  $\mathbf{p}$ 
3:  $\mathbf{r} \leftarrow \{\}$  ▷ Initialize the output string
4:  $q_{curr} \leftarrow q_0$  ▷ Initialize the current dfa state the response is at
5: for  $i \in \{1, \dots, k\}$  do
6:    $\mathbf{x} \cdot \mathbf{r}_i, q_{next} \leftarrow \text{Diffusion}(\mathbf{x}, m + (i - 1) \times d, d, T, V, Q_l, \delta, q_{curr})$ 
7:   if  $q_{next} \notin Q_l$  then ▷ No valid completion
8:     return None
9:    $\mathbf{r} \leftarrow \mathbf{r} \cdot \mathbf{r}_i$ 
10:   $\mathbf{x} \leftarrow \mathbf{x} \cdot \mathbf{r}_i$  ▷ Compute the input prompt for the next block
11:   $q_{curr} \leftarrow q_{next}$  ▷ Update current DFA state for next block
12: Return  $\mathbf{r}$ 
```

903 In the semi-autoregressive setting, after each block, we ensure that the output generated so far ends in
904 a live state from Q_l ; otherwise, we return the None string (line 7, Algo. 5). Additionally, we maintain
905 a variable q_{curr} to track the current DFA state at the end of each block. This state is then used as the
906 starting state for the dynamic programming step in the constrained generation of the next block.

907 E Token Transitions Statistics

Table 3: Token Transitions Pre-Computation Statistics

Model Family	V	GSM-Symbolic		JSON-Mode	
		Time(s)	#States	Time(s)	#States
LLaDA-8B	126349	32.09	40	13.22	169.31
Dream-7B	151667	37.01	40	11.87	169.31

908 In Table 3, we report the precomputation time and the number of states in the DFA for both tasks.
 909 For JSON generation, different regular expressions are used for different schemas; therefore, we
 910 report the mean precomputation time and mean number of states. The maximum number of states
 911 and precomputation times across all questions are 455 and 17.7 (Dream) 21.3 (LLaDA) seconds,
 912 respectively.

913 F GSM-Symbolic

914 F.1 GSM-Symbolic Prompt

```

915 You are an expert in solving grade school math tasks. You will be presented
916 with a grade-school math word problem with symbolic variables and be
917 asked to solve it.
918
919 Before answering you should reason about the problem (using the <reasoning>
920 field in the response described below). Intermediate symbolic expressions
921 generated during reasoning should be wrapped in << >>.
922
923 Only output the symbolic expression wrapped in << >> that answers the
924 question. The expression must use numbers as well as the variables
925 defined in the question. You are only allowed to use the following
926 operations: +, -, /, //, %, *, and **.
927
928 You will always respond in the format described below:
929 Let's think step by step. <reasoning> The final answer is <<symbolic
930 expression>>
931
932 There are {t} trees in the {g}. {g} workers will plant trees in the {g} today
933 . After they are done, there will be {tf} trees. How many trees did the {
934 g} workers plant today?
935
936 Let's think step by step. Initially, there are {t} trees. After planting,
937 there are {tf} trees. The number of trees planted is <<tf - t>>. The
938 final answer is <<tf - t>>.
939
940 If there are {c} cars in the parking lot and {nc} more cars arrive, how many
941 cars are in the parking lot?
942
943 Let's think step by step. Initially, there are {c} cars. {nc} more cars
944 arrive, so the total becomes <<c + nc>>. The final answer is <<c + nc>>.
945
946 {p1} had {ch1} {o1} and {p2} had {ch2} {o1}. If they ate {a} {o1}, how many
947 pieces do they have left in total?
948
949 Let's think step by step. Initially, {p1} had {ch1} {o1}, and {p2} had {ch2}
950 {o1}, making a total of <<ch1 + ch2>>. After eating {a} {o1}, the
951 remaining total is <<ch1 + ch2 - a>>. The final answer is <<ch1 + ch2 - a
952 >>.
953
954 {p1} had {l1} {o1}. {p1} gave {g} {o1} to {p2}. How many {o1} does {p1} have
955 left?
956
957 Let's think step by step. {p1} started with {l1} {o1}. After giving {g} {o1}
958 to {p2}, {p1} has <<l1 - g>> {o1} left. The final answer is <<l1 - g>>.
959
960 {question}
961

```

Listing 1: Prompt template for the GSM-Symbolic task Mirzadeh et al. [2024].

963

964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
~~994~~
~~995~~

Listing 2: GSM-Symbolic Regex

996

997

Question: {s1} gets {b} emails a day. {c}% of those emails don't require any response. He responds to the rest of them. How many emails does he respond to in a {a} day work week?

Variables are shown in blue

Unconstrained: Let's think step by step. ``Reasoning Text ...`` that require a response is $\langle\langle b \rangle\rangle * a \rangle\rangle = \langle\langle b \rangle\rangle * \langle\langle b \rangle\rangle * a \rangle\rangle$. The final answer is $\langle\langle b \rangle\rangle * \langle\langle b \rangle\rangle * a \rangle\rangle$.

Syntax error 😞

Greedy Constrained: Let's think step by step. ``Reasoning Text ...`` of emails requiring require a response is $\ll b - (c/100 * b^{**} b * a) ** a \gg$. The final answer is $\ll b - (c/100 * b^{**} b * a) ** a \gg$.

Functionally incorrect 😞

DINGO: Let's think step by step. ``Reasoning Text
...``The number of emails responded to is {b -
<<c/100>> * b} * a. The final answer is <<(b - (c/100)
* b) * a>>.

Syntactically and Functionally Correct 😊

Figure 2: An example from the GSM-symbolic dataset (variables in blue), where unconstrained generation produces syntactically incorrect output, and greedy constrained generation yields a syntactically valid but incorrect answer. In contrast, DINGO generates the correct answer.

998

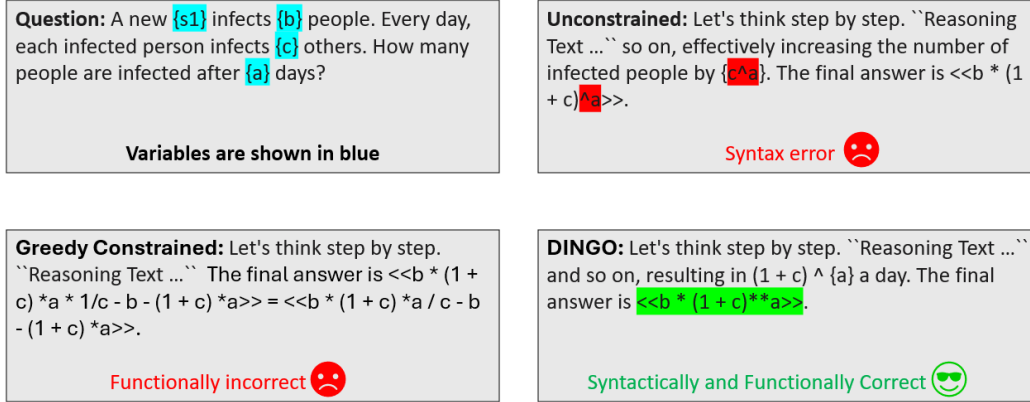


Figure 3: An example from the GSM-symbolic dataset (variables in blue), where unconstrained generation produces syntactically incorrect output, and greedy constrained generation yields a syntactically valid but incorrect answer. In contrast, DINGO generates the correct answer.

G JSON-Mode

G.1 JSON-Mode Example Prompt

```

1001 You are a helpful assistant that answers in JSON. Here's the json schema you
1002 must adhere to:
1003 <schema>
1004 {
1005   'title': 'PromotionalCampaign', 'type': 'object', 'properties': {
1006     'campaignID': {
1007       'title': 'Campaign ID', 'type': 'string', 'productID': {
1008         'title': 'Product ID', 'type': 'string', 'startDate': {
1009           'title': 'Start Date', 'type': 'string', 'format': 'date', 'endDate': {
1010             'title': 'End Date', 'type': 'string', 'format': 'date', 'discountDetails': {
1011               'title': 'Discount Details', 'type': 'string'
1012             }
1013           }, 'required': ['campaignID', 'productID', 'startDate', 'endDate']
1014         }
1015       }
1016     }
1017   }
1018 }
1019 </schema>
1020 I'm organizing a promotional campaign for our new eco-friendly laundry
1021 detergent, which is part of our household products line. The campaign
1022 will start on June 1, 2023, and end on June 30, 2023. We're planning to
1023 offer a 15% discount on all purchases during this period. The campaign ID
1024 is CAMP123456, the product ID is PROD7891011, and the discount details
1025 are 15% off on all purchases.
1026 Only output the JSON object, no other text or comments.

```

Listing 3: Example JSON Prompt from the JSON-Mode-Eval task NousResearch [2024]. The prompt includes a system message that specifies a schema and a user message that explicitly instructs the model to output a JSON object following that schema with certain parameters.

G.2 JSON-Mode Example Regex

```

1023 \\{[ ]?"campaignID"[ ]?:[ ]?"([~\\\\\\\\\\\\\\\\x00-\\\\x1F\\\\x7F-\\\\x9F]|\\\\\\\\[~\\\\\\\\\\\\\\\\])*"
1024 [ ]?,[ ]?"productID"[ ]?:[ ]?"([~\\\\\\\\\\\\\\\\\\\\x00-\\\\x1F\\\\x7F-\\\\x9F]|\\\\\\\\[~\\\\\\\\\\\\\\\\])
1025 *"[ ]?,[ ]?"startDate"[ ]?:[ ]?"(?:\\\\d{4})-(?:0[1-9]|1[0-2])
1026 -(?:0[1-9]|1[1-2][0-9]|3[0-1])"[ ]?,[ ]?"endDate"[ ]?:[ ]?"(?:\\\\d{4})
1027 -(?:0[1-9]|1[0-2])-(?:0[1-9]|1[1-2][0-9]|3[0-1])"([ ]?,[ ]?"
1028 discountDetails"[ ]?:[ ]?"([~\\\\\\\\\\\\\\\\\\\\x00-\\\\x1F\\\\x7F-\\\\x9F]|\\\\\\\\[~\\\\\\\\\\\\\\\\])*"
1029 [ ]?\\}
1030

```

Listing 4: Regex for the JSON Schema in Appendix G.2

G.3 JSON-Mode Case Studies

Question: Here's the json schema you must adhere to: `<schema> ... </schema>` <Question ... Please respond with a valid JSON object for this order.> Only output the JSON object, no other text or comments.

```
Unconstrained: {'orderId': 'ORD10234', 'items':
  {'name': 'Wireless Mouse', 'price': 15.99}, {'name':
  'Bluetooth Keyboard', 'price': 29.99}, {'name':
  'Webcam HD', 'price': 45.50}}
```

Syntax error (missing ']') 😞

Greedy Constrained: {"orderId": "4", "items": [{"name": "Wireless Mouse", "price": 15.99}, {"name": "Bluetooth Keyboard", "price": 29.99}, {"name": "Webcam HD", "price": 45.50}]} 1

Syntax error (incomplete output) 😞

DINGO: {'orderId': 'ORD10234', 'items': [{'name': 'Wireless Mouse', 'price': 15.99}, {'name': 'Bluetooth Keyboard', 'price': 29.99}, {'name': 'Webcam HD', 'price': 45.50}]}

Syntactically Correct 😎

Figure 4: An example from JSON generation, where unconstrained generation produces a syntactically incorrect output, and greedy constrained generation yields a valid but incomplete prefix. In contrast, DINGO generates a syntactically correct answer.

Question: Here's the json schema you must adhere to:<schema> ... </schema> Question ... service was very prompt and professional, but I wish there was more variety in the product options.'Only output the JSON object, no other text or comments.

Unconstrained: `\n\n"customerID": "CUST29382",
"feedbackScore": 4, "comments": "The service was
very prompt and professional, but I wish there was
more variety in the product options."}`

Syntax error (missing '{') 😞

Greedy Constrained: { "customerID":
"1UST9382", "feedbackScore": 4,
\n\n\n\n\n"comments"\n\n\n\n\n\n\n". Think
like the are D to a five year Bible. ... a little about the
and the human world. Is the author review

Syntax error (incomplete output) 😞

DINGO: {'customerID': 'CUST29382',
'feedbackScore': 4, 'comments': 'The service was very
prompt and professional, but I wish there was more
variety in the product options.'}

Syntactically Correct 😎

Figure 5: An example from JSON generation, where unconstrained generation produces a syntactically incorrect output, and greedy constrained generation yields a valid but incomplete prefix. In contrast, DINGO generates a syntactically correct answer.

1033 **H Ablation Study on Number of Blocks for Diffusion LLM Generation**
1034 **(GSM-Symbolic)**

1035 We run generation with a response length of 128, using 64 total diffusion steps, and each of 1, 2, and
1036 8 blocks. Table 4 presents the result.

Table 4: Ablation Study on The Number of Diffusion Blocks for GSM-Symbolic

Model	#Blocks	Method	Acc. (%)	Parse (%)	Time (s)
LLaDA-8B-I	1	Unconstrained	20	54	23.66
		Greedy Constrained	26	94	23.7
		Best of Greedy + Unconstrained	26	94	23.66
		DINGO	29	100	23.73
	2	Unconstrained	22	54	23.63
		Greedy Constrained	30	96	23.81
		Best of Greedy + Unconstrained	30	96	23.65
		DINGO	32	100	23.93
	8	Unconstrained	19	35	23.78
		Greedy Constrained	27	98	23.97
		Best of Greedy + Unconstrained	27	98	23.8
		DINGO	32	100	23.92
	1	Unconstrained	28	69	23.56
		Greedy Constrained	32	90	23.64
		Best of Greedy + Unconstrained	32	90	23.65
		DINGO	34	100	23.67
Dream-I-7B	2	Unconstrained	30	55	23.62
		Greedy Constrained	33	87	23.71
		Best of Greedy + Unconstrained	33	87	23.62
		DINGO	34	100	23.65
	8	Unconstrained	32	61	23.89
		Greedy Constrained	34	93	24.01
		Best of Greedy + Unconstrained	34	93	23.89
		DINGO	36	100	23.91

1037 **I Ablation Study on Number of Blocks for Diffusion LLM Generation**
1038 **(JSON-Mode)**

1039 We run generation with a response length of 128, using 64 total diffusion steps, and each of 1, 2, and
1040 8 blocks. Table 5 presents the result.

Table 5: Ablation Study on The Number of Diffusion Blocks for JSON-Mode.

Model	#Blocks	Method	Acc. (%)	Parse (%)	Time (s)
LLaDA-8B-I	1	Unconstrained	87	91	6.7
		Greedy Constrained	78	79	6.81
		Best of Greedy + Unconstrained	99	99	6.73
		DINGO	100	100	6.78
	2	Unconstrained	84	92	6.72
		Greedy Constrained	92	94	6.83
		Best of Greedy + Unconstrained	99	99	6.73
		DINGO	100	100	6.86
	8	Unconstrained	84	89	6.73
		Greedy Constrained	98	98	6.87
		Best of Greedy + Unconstrained	100	100	6.75
		DINGO	100	100	6.85
Dream-I-7B	1	Unconstrained	85	87	6.4
		Greedy Constrained	30	30	6.51
		Best of Greedy + Unconstrained	91	93	6.43
		DINGO	100	100	6.55
	2	Unconstrained	79	82	6.47
		Greedy Constrained	37	39	6.68
		Best of Greedy + Unconstrained	86	88	6.5
		DINGO	100	100	6.63
	8	Unconstrained	70	74	6.44
		Greedy Constrained	52	52	6.65
		Best of Greedy + Unconstrained	86	89	6.46
		DINGO	100	100	6.67

1041 **J Ablation Study on Number of Steps for Diffusion LLM Generation**
1042 **(GSM-Symbolic)**

1043 We run generation with a response length of 128, 1 block, and each of 16, 32, 64, and 128 total
1044 diffusion steps. Table 6 presents the result.

Table 6: Ablation Study on The Number of Diffusion Steps for GSM-Symbolic with Dream-I-7B

#Steps	Method	Acc. (%)	Parse (%)	Time (s)
16	Unconstrained	6	20	5.99
	Greedy Constrained	13	78	6.18
	Best of Greedy + Unconstrained	13	78	5.99
	DINGO	18	100	6.09
32	Unconstrained	18	48	11.96
	Greedy Constrained	25	87	12.06
	Best of Greedy + Unconstrained	25	87	11.96
	DINGO	28	100	12.03
64	Unconstrained	28	69	23.56
	Greedy Constrained	32	90	23.64
	Best of Greedy + Unconstrained	32	90	23.65
	DINGO	34	100	23.67
128	Unconstrained	31	74	47.83
	Greedy Constrained	30	89	47.88
	Best of Greedy + Unconstrained	31	90	47.83
	DINGO	33	100	47.86

1045 **K Ablation Study on Number of Steps for Diffusion LLM Generation**
1046 **(JSON-Mode)**

1047 We run generation with a response length of 128, 1 block, and each of 16, 32, 64, and 128 total
1048 diffusion steps. Table 7 presents the result.

Table 7: Ablation Study on The Number of Diffusion Steps for JSON-Mode with Dream-I-7B

#Steps	Method	Acc. (%)	Parse (%)	Time (s)
16	Unconstrained	54	59	1.51
	Greedy Constrained	32	32	1.62
	Best of Greedy + Unconstrained	68	71	1.52
	DINGO	100	100	1.6
32	Unconstrained	67	71	3.23
	Greedy Constrained	35	35	3.35
	Best of Greedy + Unconstrained	78	82	3.24
	DINGO	100	100	3.31
64	Unconstrained	85	87	6.4
	Greedy Constrained	30	30	6.51
	Best of Greedy + Unconstrained	91	93	6.43
	DINGO	100	100	6.55
128	Unconstrained	85	87	13.42
	Greedy Constrained	46	46	13.53
	Best of Greedy + Unconstrained	95	97	13.43
	DINGO	100	100	13.51